

Bioinformatique: prédiction de gènes

INSA - Février 2004

Laurent Duret

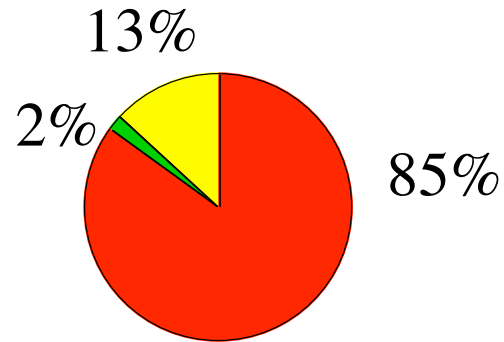
BBE – UMR CNRS n° 5558

Université Claude Bernard - Lyon 1

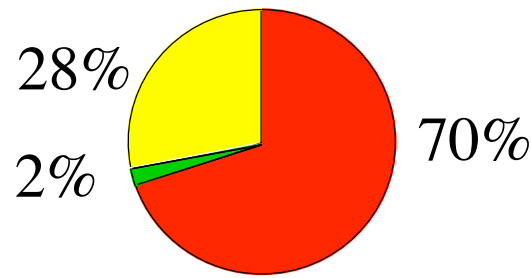
Genome annotation

- √ Identification of repeats (RepeatMasker, Reputer, ...)
- √ Prediction of protein-coding genes
 - λ Intrinsic methods (GenScan, Genmark, Glimmer, ...)
 - λ Genomic/mRNA (EST) comparison (blastn, sim4, ...)
 - λ Genomic/protein comparison (blastx, GeneWise, ...)
- √ Prediction of RNA genes
 - λ Intrinsic methods (tRNA: tRNAScanSE, snoRNA ...)
 - λ Genomic/RNA (EST) comparison (blastn, sim4, ...)
- √ And more ...
 - λ Replication origins (bacteria) (oriloc)
 - λ Pseudogenes (by similarity) (blastn, blastx)
 - λ Regulatory elements (CpG islands, promoters ??)

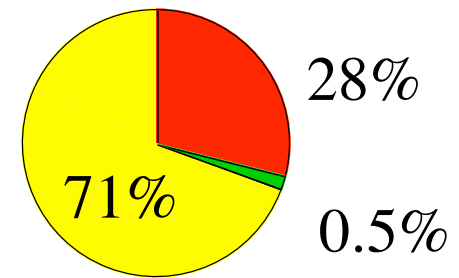
Proportion of functional elements within genomes



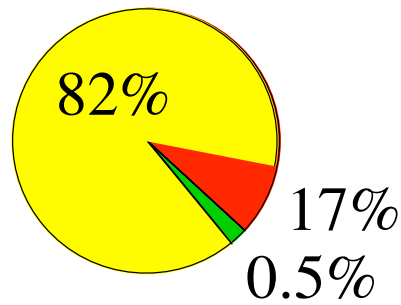
E. coli



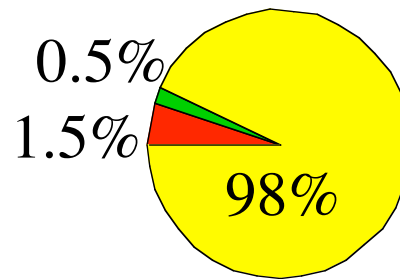
Yeast
S. cerevisiae



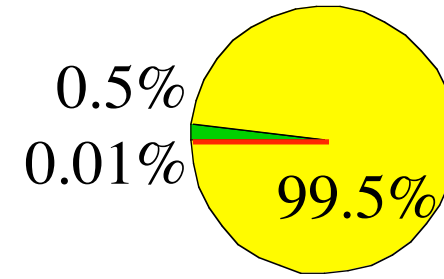
Nematode
C. elegans



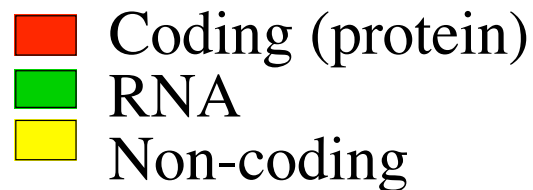
Drosophila



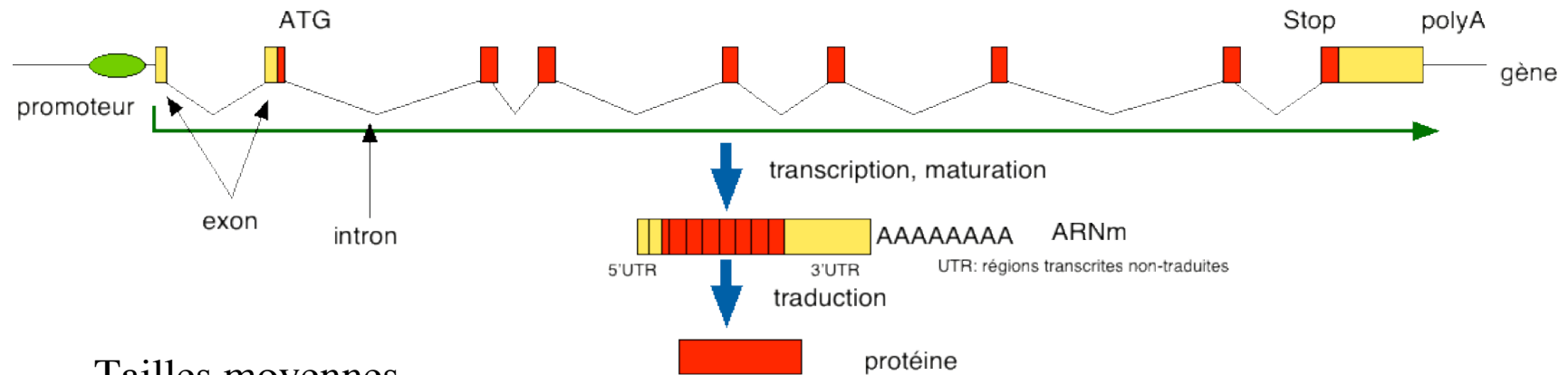
Human



Lunfish
(dipnoi)



Structure des gènes humains



v Tailles moyennes

λ Gene	27 kb
λ CDS	1100 nt
λ Exon (interne)	145 nt
λ Intron	3400 nt
λ 5'UTR	210 nt
λ 3'UTR	740 nt

v Intron/exon

λ Nombres d'introns:	6 ± 3 introns / kb CDS
λ Introns / (introns + CDS):	80%

v Epissage alternatif dans plus de 30% des gènes

Prédiction de gènes: informations utilisées

- ✓ 1- caractérisation de la taille et du contenu des régions (codantes/non-codantes)
 - ✓ 2- caractérisation des signaux au niveau de sites fonctionnels (e.g. signaux d'épissage, début et fin de traduction, ...)
 - ✓ 3- utilisation de similarité ADN/protéines, ADN/ARNm, ADN/ADN
-
- ✓ méthodes intrinsèques (ab initio): utilisent 1 et 2
 - ✓ méthodes extrinsèques (approche comparative): utilisent 3, et éventuellement 2

Prédiction de gènes : méthodes intrinsèques

- √ Prédiction des régions codantes uniquement !
- √ Recherche de phases ouvertes de lecture (ORF: open reading frame) = série de codon sans STOP

Phase +0

Phase +1

Phase +2

ATGTACCGTCGATCGTAGCTTGATCGATCG

TACATGGCAGCTAGCATCGAACTAGCTAGC

Phase -0

Phase -1

Phase -2

- Taille moyenne des ORF: ± 150 nt

- √ Distinction codant/non-codant : contenu et taille des séquences
 - λ usage des codons: utilisation non aléatoire des codons synonymes
 - λ fréquence des amino-acides (e.g. tryptophane est rare)
 - λ corrélations entre amino-acides (codons) successifs
 - λ taille des exons et introns

- λ Apprentissage sur un ensemble de gènes connus
- λ Fréquence d'oligomères (e.g. hexamères)
- λ chaînes de Markov

Prédiction de gènes : méthodes intrinsèques (suite)

- √ Recherche de signaux: sites fonctionnels conservés
 - λ signaux d'épissage: site donneur, accepteur d'épissage, point de branchement
 - λ codon d'initiation de la traduction
 - λ codon stop

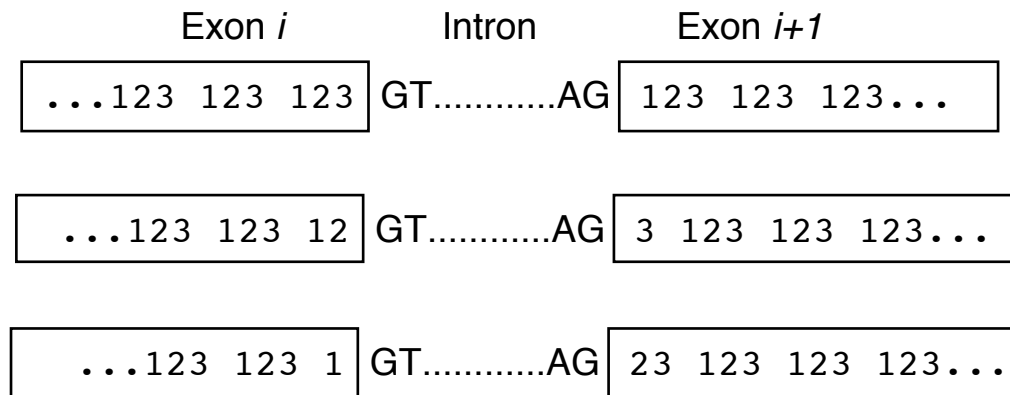
- λ Utilisation de consensus (historique): e.g.

donneur	accepteur
A/CAG GT RAGT	YYYYYYYYYY*C AG G

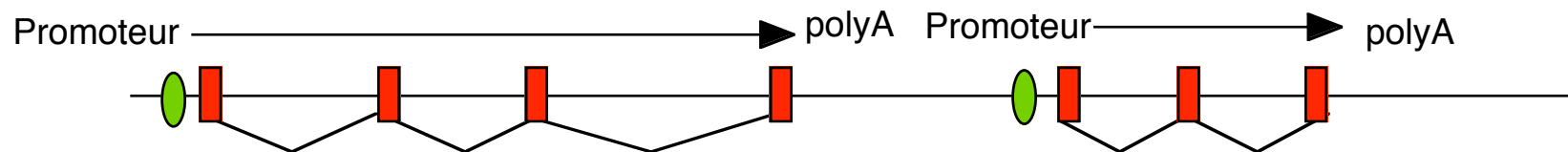
- λ Utilisation de matrices de pondération position-dépendantes (profils): Cf suite du cours

Prédiction de gènes : méthodes intrinsèques (suite)

- Construction d'un modèle de gène protéique
 - Combinaison d'exons de phases compatibles (pondération en fonction des scores de chaque exon potentiel) - pas de codons stop en phase!



- Recherche de limites de gènes
 - Exons terminaux (5', 3')
 - Promoteur
 - Signal de polyadénylation



Qualité de la prédiction par exon

- ✓ Évaluation de la fiabilité de la prédiction
 - λ essai des logiciels de prédiction sur un ensemble de séquences caractérisées expérimentalement (différentes de celles utilisées pour entraîner les logiciels)

- ✓ Sensibilité : fraction des exons présents dans la séquence qui sont retrouvés par le logiciel

$$\text{sensibilité} = \frac{\text{Vrais Positifs}}{\text{Vrais Positifs} + \text{Faux Négatifs}}$$

- λ e.g. GenScan (mammifères): 78%

- ✓ Spécificité : fraction des vrais exons parmi tous ceux prédits

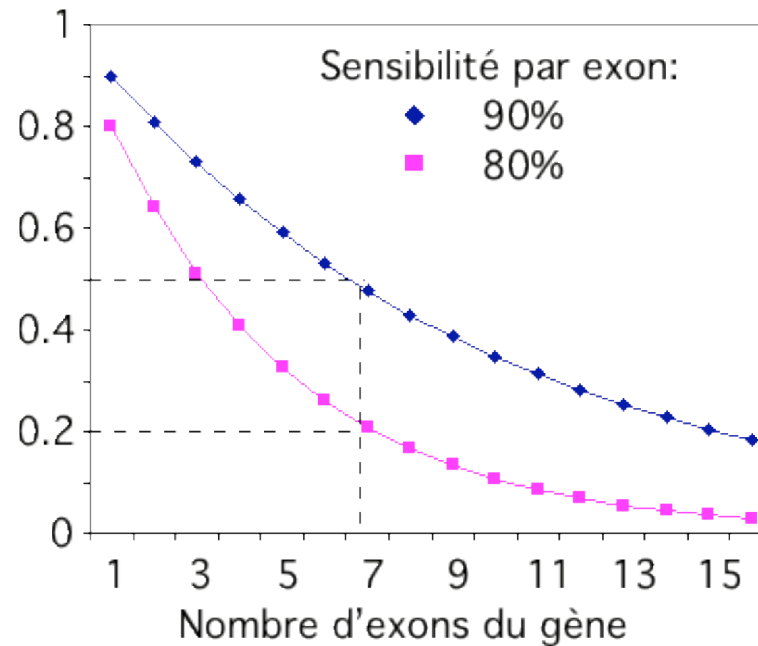
$$\text{spécificité} = \frac{\text{Vrais Positifs}}{\text{Vrais Positifs} + \text{Faux Positifs}}$$

- λ e.g. GenScan (mammifères): 81%

Prédiction de gènes protéiques complets

- Construction d'un modèle de gène à partir de prédictions d'exons de phases compatibles
- Prédiction de gènes complets: sensibilité ?

Probabilité de détecter tous les exons d'un gènes



λ + les faux positifs ! + épissage alternatif ! + exons non-codants !

Un peu d'optimisme

- ∨ Fraction de la longueur des gènes correctement prédits:

70-80%

- ∨ Probabilité que deux exons potentiels consécutifs soient réels (et donc positifs en RT-PCR)

0.5

Prédiction de gènes : méthodes intrinsèques (bilan)

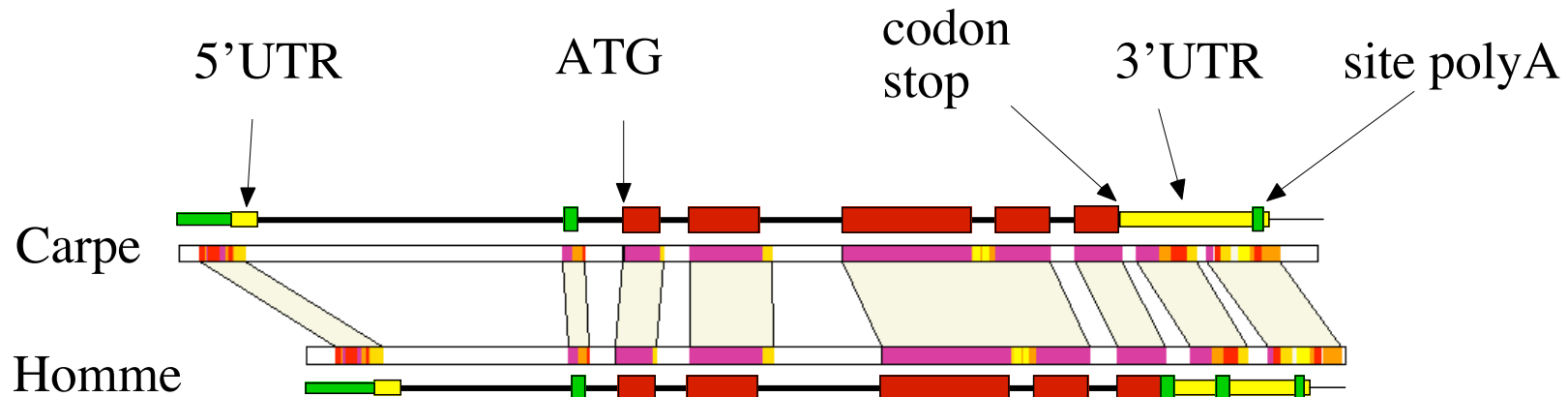
- √ Procaryotes (pas d'intron):
 - λ sensibilité et spécificité > 95% (dépend du taux de G+C du génome)
- √ Eucaryotes: efficacité variable (dépend du taux de G+C du génome et du nombre et de la taille des introns)
 - λ prédiction d'exons: sensibilité et spécificité 60-80%
 - λ prédiction de gènes complets:
 - levure: >90% des gènes correctement prédits
 - nématode: 50% des gènes correctement prédits
 - homme: 20% (?) des gènes correctement prédits
- √ très utile pour guider les expérimentations

Prédiction de gènes : méthodes extrinsèques

- √ Utilisation des EST
 - λ comparaison séquence ADN génomique / mRNA : identification des exons (blastn, sim4)
 - λ informations sur épissage alternatif, expression
 - λ problème:
 - gènes faiblement exprimés ou à distribution tissulaire restreinte
 - artéfacts dans les EST

- √ Approche comparative
 - λ Comparaison d'une séquence génomique avec des gènes déjà caractérisés dans d'autres espèces (ADN/protéine) (blastx, genewise)
 - λ Comparaison de séquences génomiques homologues (ADN/ADN)

Analyse comparative des gènes de β -actine de l'homme et de la carpe



introns: —
régions codantes: ■
éléments régulateurs: ■

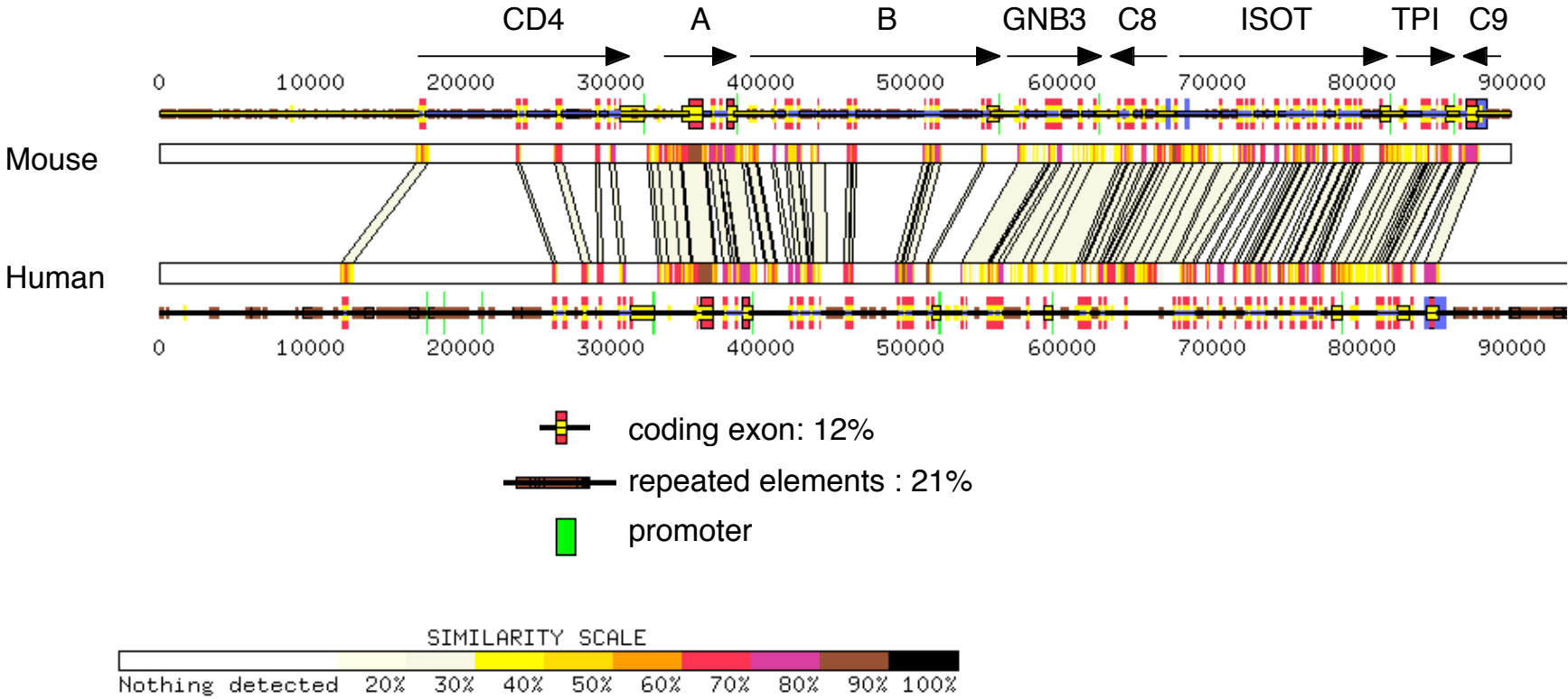
échelle de similarité:

□ pas de similarité significative
■ 80 - 90% identité
■ 70 - 80% identité

Comparison of human and mouse CD4-C9 locus:
gene-rich, repeated-element poor, G+C-rich region (50.5%)

Human chromosome 12p13
Mouse chromosome 6

8 genes: CD4, A, B, GNB3, C8, ISOT, TPI, C9



Prédiction de gènes : démarche

- √ 1- recherche de séquences répétées (RepeatMasker)
- √ 2- méthodes intrinsèques (consensus de différentes méthodes)
- √ 3- recherche de similarité ADN/protéines (blastx/genewise)
- √ 4- recherche de similarité ADN/mRNA (blastn/sim4)
- √ 5- recherche de similarité ADN/ADN (blastn)
- √ COMBINER LES RESULTATS

- √ 6- prédiction de gènes RNA
 - λ tRNA: tRNAScanSE
 - λ rRNA: par similarité
 - λ snRNA ...

Prédiction de régions régulatrices

- √ Méthodes intrinsèques (*ab initio*)

- λ Prédiction de promoteurs

- λ Îlots CpG

- √ Approche comparative

Prédiction de promoteurs eucaryotes

- √ Combinaison de sites de fixation de facteur de transcription (ordre, orientation, distance)
- √ Motifs courts, dégénérés
 - λ Difficile de distinguer les vrais sites des faux positifs:
 - λ Motif à 4 bases: $\approx 1/256$ pb (1/128 pb sur les deux brins)
- √ Boîtes TATA, CAAT , GC: absents dans beaucoup de promoteurs

- √ Banques de données de sites de fixation de facteurs de transcription (TRANSFAC), de promoteurs caractérisés expérimentalement (EPD)

- √ PromoterScan (Prestridge 1995): Mesure de la densité en sites potentiels de fixation de facteurs de transcription de long de la séquence (pondération en fonction de la fréquence des sites dans ou en dehors des vrais promoteurs)

Prédiction de promoteurs: sensibilité, spécificité

- v Sensibilité: fraction des promoteurs qui sont trouvés par le logiciel

$$sensibilité = \frac{vrais_positifs}{vrais_positifs + faux_négatifs}$$

- λ PromoterScan: sensibilité = 70% (promoteurs à boîte TATA)

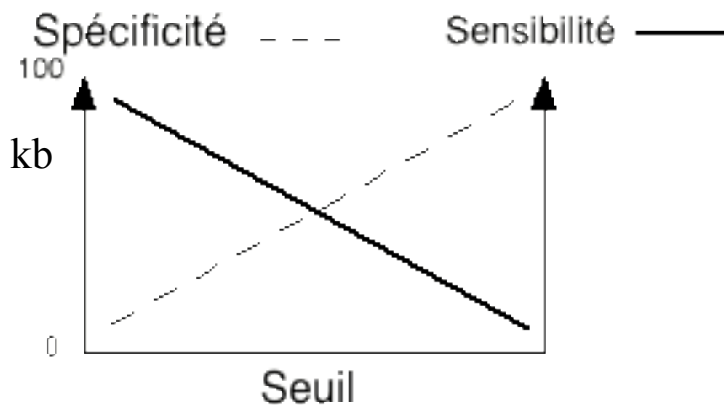
- v Spécificité: fraction des vrais promoteurs parmi ceux qui ont été prédits

$$spécificité = \frac{vrais_positifs}{vrais_positifs + faux_positifs}$$

- λ PromoterScan: spécificité = 20%

- λ Un faux positif / 10 kb

- v Génome humain: $\approx 30\,000$ gènes, ≈ 1 promoteur/100 kb



Prédiction de promoteurs eucaryotes: recherches en cours

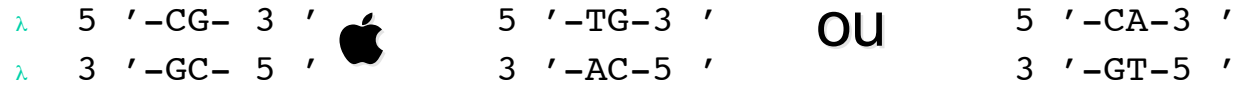
- v Prise en compte de l'orientation relative et des distances entre sites de fixation de facteurs de transcription
 - λ COMPEL (Kolchanov 1998): banque de données d'éléments composites
 - λ FastM : recherche dans une séquence génomique d'une combinaison de deux sites de fixation de facteurs de transcription à une distance définie l'un de l'autre

 - v Recherche de corrélations entre sites
 - λ PromoterInspector (Werner 2000)
 - Sensibilité: 40%
 - Spécificité: 45%
- <http://www.gsf.de/biodv/index.html>
-
- v Combinaison recherche *ab initio* / approche comparative: recherche de sites potentiels parmi les régions conservées

Îlots CpG

- Genome de vertébrés :
 - λ méthylation des C dans les dinucléotides 5'-CG-3' (CpG)

- Me-C fortement mutable -> T



- Genome des vertébrés: globalement dépourvu en CpG (excès de TG, CA)

$$CpG_{o/e} = \frac{\text{Nombre_de_CpG_observé}}{\text{Nombre_de_CpG_attendu}} = 0.25$$

- Certaines régions (200 nt à plusieurs kb) échappent à la méthylation

- λ Pas de déplétion en CpG: CpG_{o/e} proche de 1

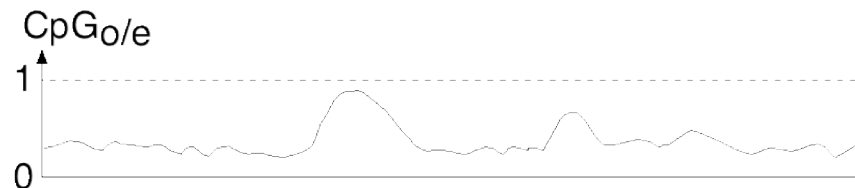
- λ Riche en G+C

- λ Îlot CpG:

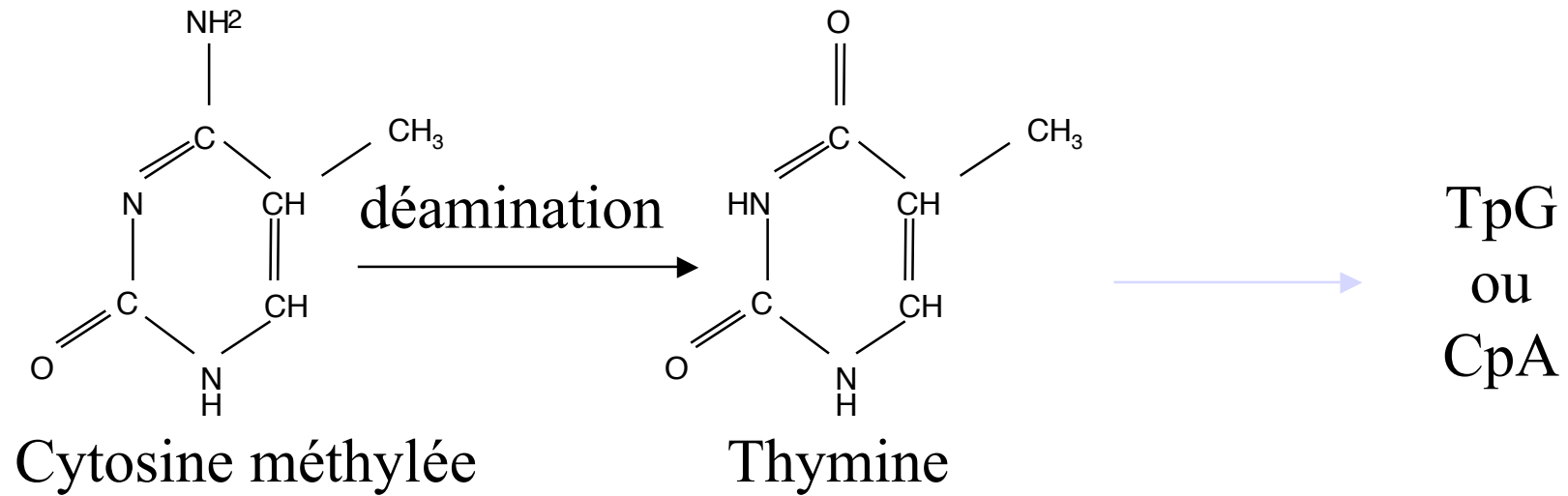
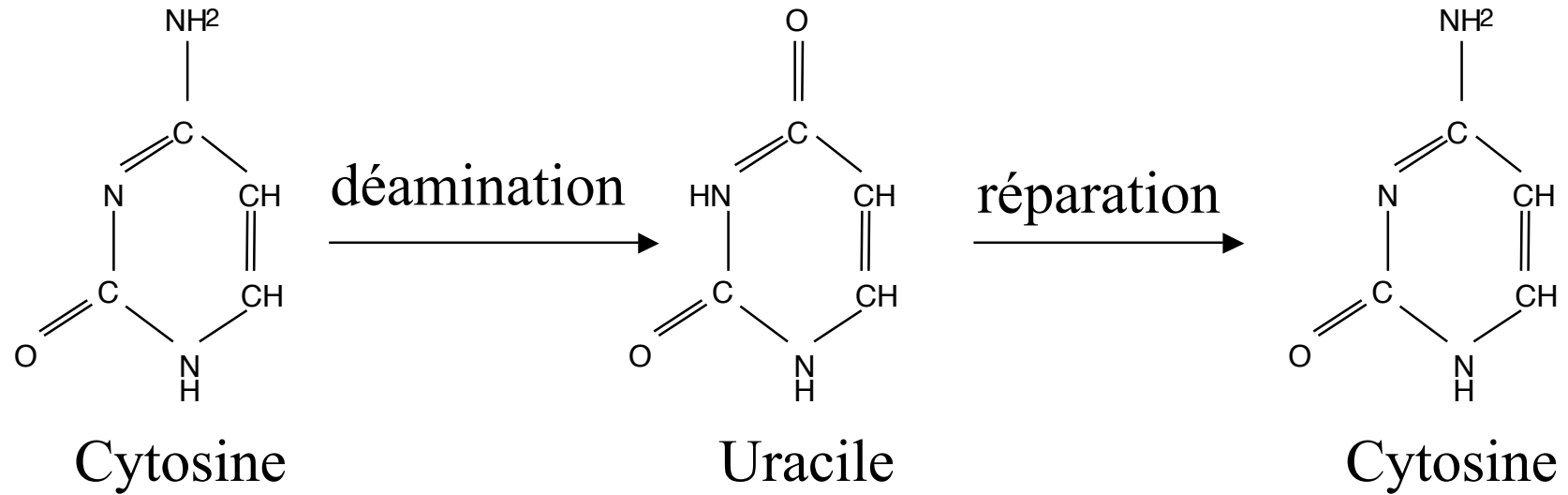
Longueur > 500 nt

CpG_{o/e} > 0.6

G+C > 50%



La déamination des cytosines



Îlots CpG: associés aux régions promotrices ?

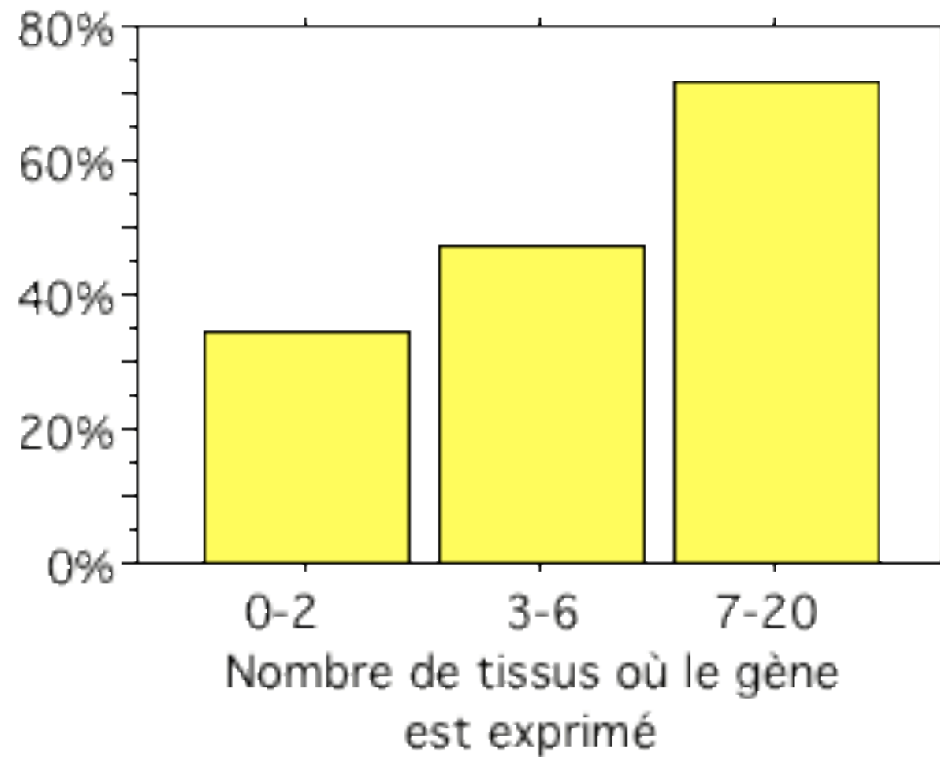
- ✓ Bird (1986), Gardiner-Garden (1987) Larsen (1992) ref
 - λ 40% des gènes tissu-spécifiques possèdent un îlot CpG en 5 '
 - λ 100% des gènes ' housekeeping ' possèdent un îlot CpG en 5 '

- ✓ Rechercher des îlots CpG pour prédire des régions promotrices ?
 - λ Sensibilité: 40-100%
 - λ Spécificité ?? (Quelle fraction des îlots CpG correspond effectivement à des régions promotrices ?)

- ✓ Ponger (2001): comparaison des îlot CpG qui recouvre ou non le site d 'initiation de la transcription

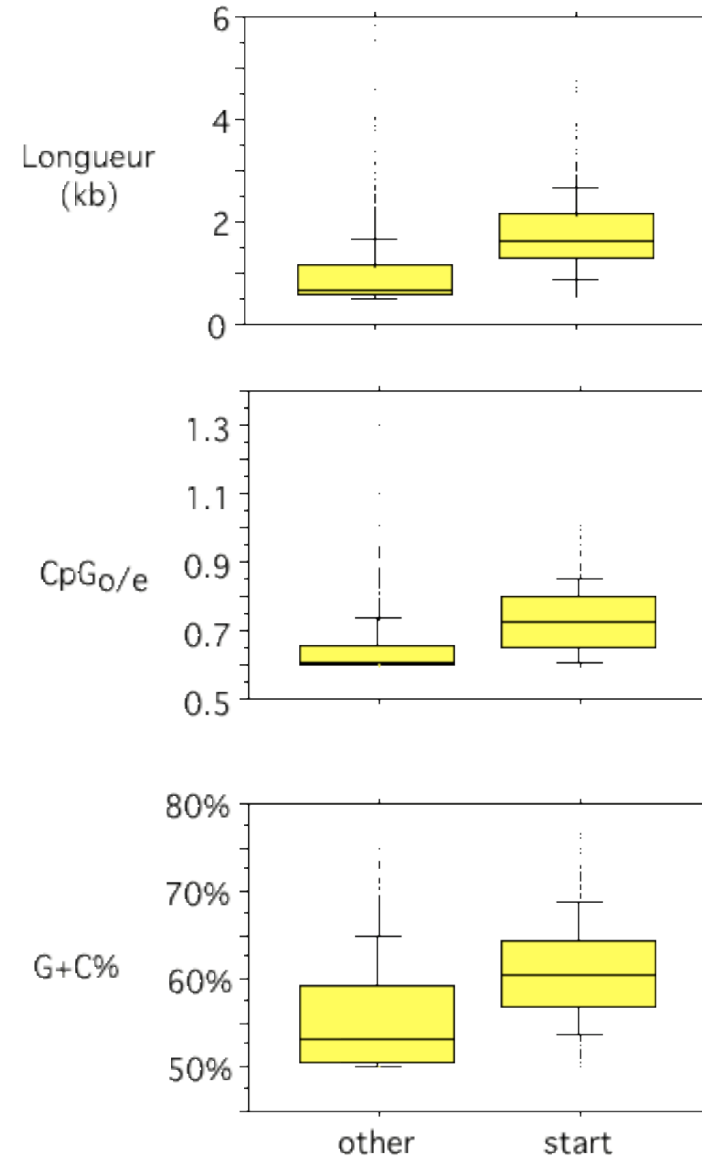
Fréquence des gènes humains avec un îlot CpG recouvrant le site d'initiation de la transcription

- λ 800 gènes humains avec promoteur décrit
- λ Mesure de la distribution tissulaire à l'aide d'EST (20 tissus)



Comparaison des îlots CpG recouvrant ou non le site d'initiation de la transcription

- λ 272 îlots start CpG recouvrant le site d'initiation de la transcription (start)
- λ 1078 îlots CpG en dehors d'un promoteur connu (other) (en excluant les séquences répétées)






Recherche de régions régulatrices par analyse comparative (empreintes phylogénétiques)

- ✓ Goodman et al. 1988: régulation de l'expression des gènes du cluster β -globine au cours du développement
 - Alignement de séquences orthologues de 6 mammifères (> 270 Ma d'évolution)
 - 13 empreintes phylogénétiques: ≥ 6 nt, conservation 100%
 - Analyse par retard de bande sur gel:
 - 12/13 (92%) correspondent à des sites de fixation de protéines
- ✓ 1996: 35 empreintes phylogénétiques avec protéines fixatrices identifiées
- ✓ Enhancers de gènes HOX (Fugu/souris) (Aparicio et al. 1995)
- ✓ enhancer TCR α (homme/souris) (Luo, 1998)
- ✓ promoteur COX5B (11 primates) (Bachman, 1996)
- ✓ promoteur uPAR (homme/souris) (Soravia, 1995)

Prediction of gene function

- √ Analysis of expression pattern (ESTs, ...)
- √ Prediction of the subcellular location of the protein : nucleus, membrane, excreted, etc.
 - λ SignalPep : <http://www.cbs.dtu.dk/services/SignalP/>
 - λ Psort: <http://psort.nibb.ac.jp/>
 - λ etc. (see <http://www.expasy.org/tools/>)
- √ Search for functional motifs (e.g. DNA binding domains, catalytic sites, ...)
 - <http://hits.isb-sib.ch/cgi-bin/PFSCAN>
- √ Prediction by homology

Function prediction by homology ?

- ✓ Similarity between proteins  homology
- ✓ Homology  conserved structure
- ✓ Conserved structure  conserved function
- ✓ Yes, but ...
 - λ Function: fuzzy concept
 - Identical biochemical activity ?
 - Identical expression pattern (tissue-specific isoforms) ?
 - Identical subcellular location (cytoplasm, mitochondria, etc.) ?
 - λ Homologous proteins with different function
 - e.g. homologous proteins binding a same receptor but opposite activity (activator/repressor)
 - homologous proteins with totally different functions: τ -crystalline / α -énolase
 - λ Orthology/paralogy
 - λ Modular evolution

Function prediction by homology ?

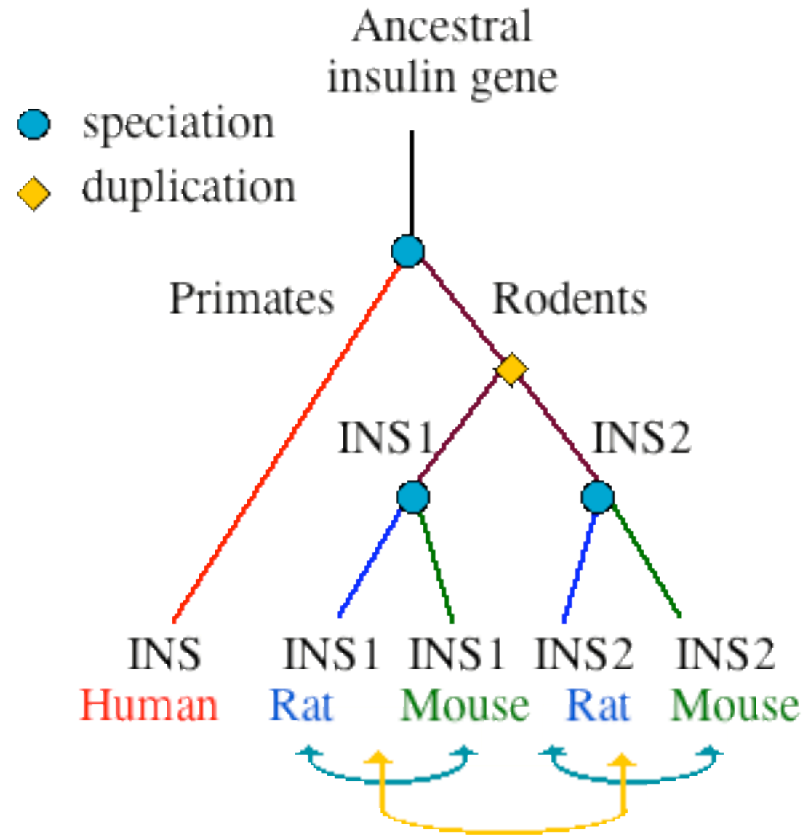
```
MZEORFG: 1   NSPDRACNLAKQAFDEAISELDSLGEESYKDSTLIMQLLDNLTTLWTSDNE 59
           N+P++AC LAKQAFD+AI+ELD+L E+SYKDSTLIMQLLDNLTTLWTSD E
BOV1433P: 186 NAPEQACLLAKQAFDDAIAELDTLNEDSYKDSTLIMQLLDNLTTLWTSEGE 244
```

```
Score = 87.4 bits (213), Expect = 1e-17
Identities = 41/59 (69%), Positives = 50/59 (84%)
```

```
LOCUS      BOV1433P      1696 bp      mRNA                      MAM      26-APR-1993
DEFINITION Bovine brain-specific 14-3-3 protein eta chain mRNA
ACCESSION  J03868
```

```
LOCUS      MZEORFG      187 bp      mRNA                      PLN      31-MAY-1994
DEFINITION Zea mays putative brain specific 14-3-3 protein,
           tau protein homolog mRNA, partial cds.
```

Orthology/paralogy



Homology: two genes are homologous if they share a common ancestor

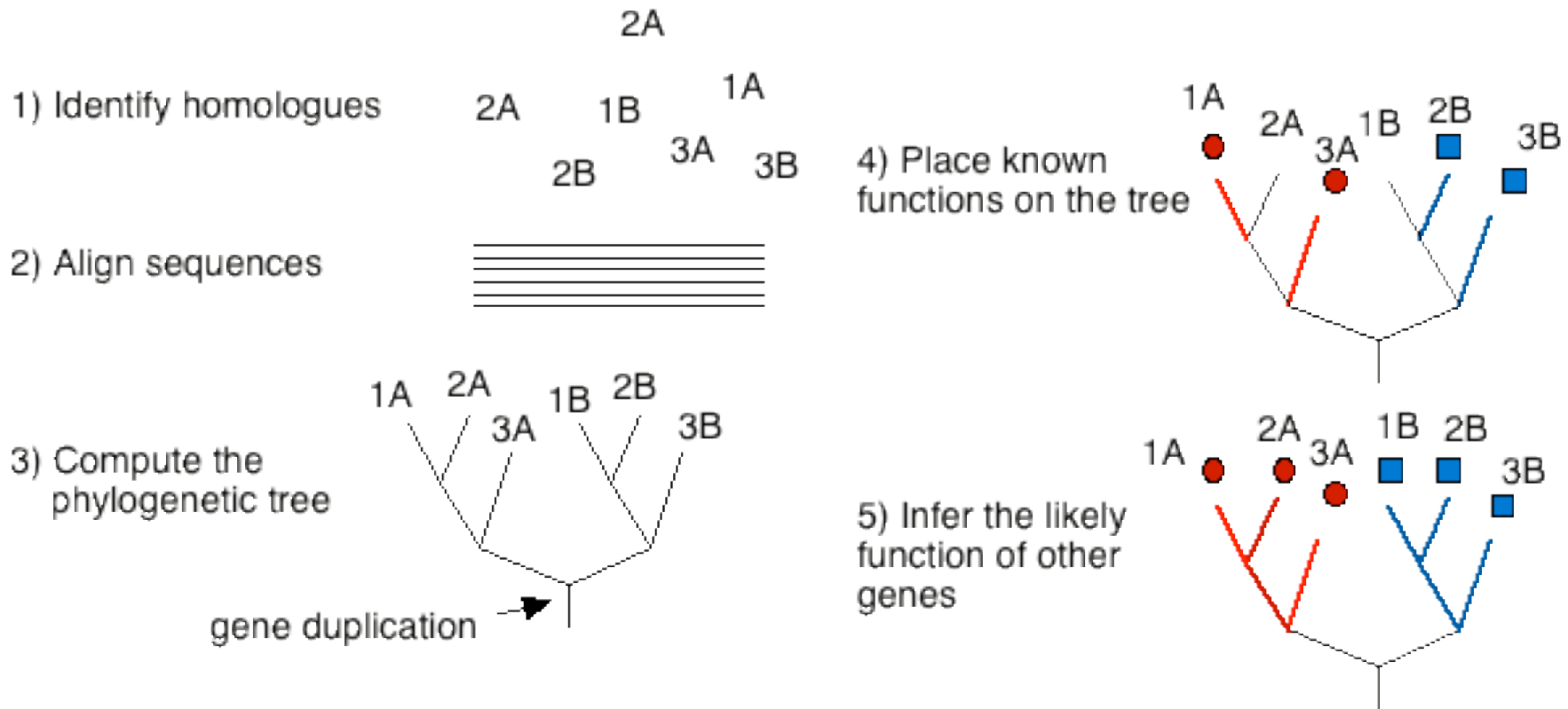
↔ Orthologues: homologous genes that have diverged after a speciation

↔ Paralogues: homologous genes that have diverged after a duplication



Orthology \neq functional equivalence

Phylogenetic approach for function prediction



Modular evolution

