

# Organisation physique et génétique des génomes

# Table des matières

- Organisation physique de l'ADN
  - Nucléosomes
  - Chromatine
- Principes d'organisation des génomes
  - Exemple d'un génome bactérien
  - Exemple du génome humain
- Evolution des génomes: duplication et recombinaison

# Chromosomes

- Structures organisées qui contiennent l'ADN

Les molécules d'AND sont arrangées avec des protéines appelées histones de manière à tenir dans un espace très réduit. L'ensemble ADN + histones est appelé chromatine

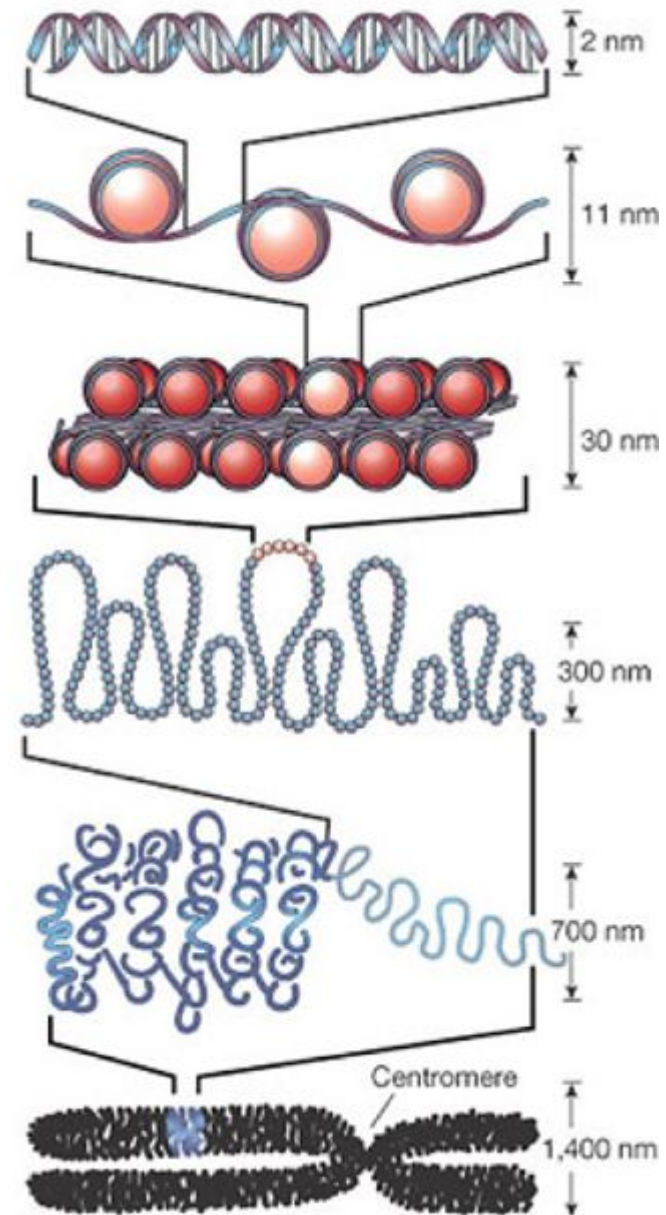
- On peut avoir plusieurs chromosomes dans le noyau

Ils sont organisés en paires de copies parentales (diploïde) ou seuls (haploïde) ou encore en nombre pair de copies quelconque (polyploïde)

# La structure du chromosome eucaryote

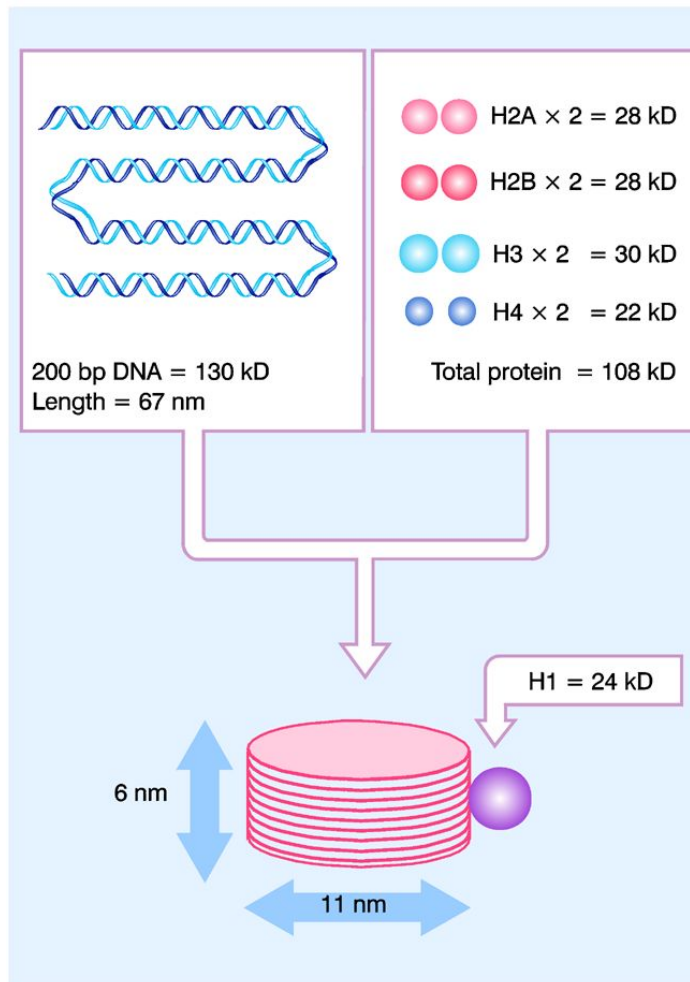
Différents niveaux de compaction

Chromatine =  
ADN  
+ protéines  
(histones +  
non-histones)



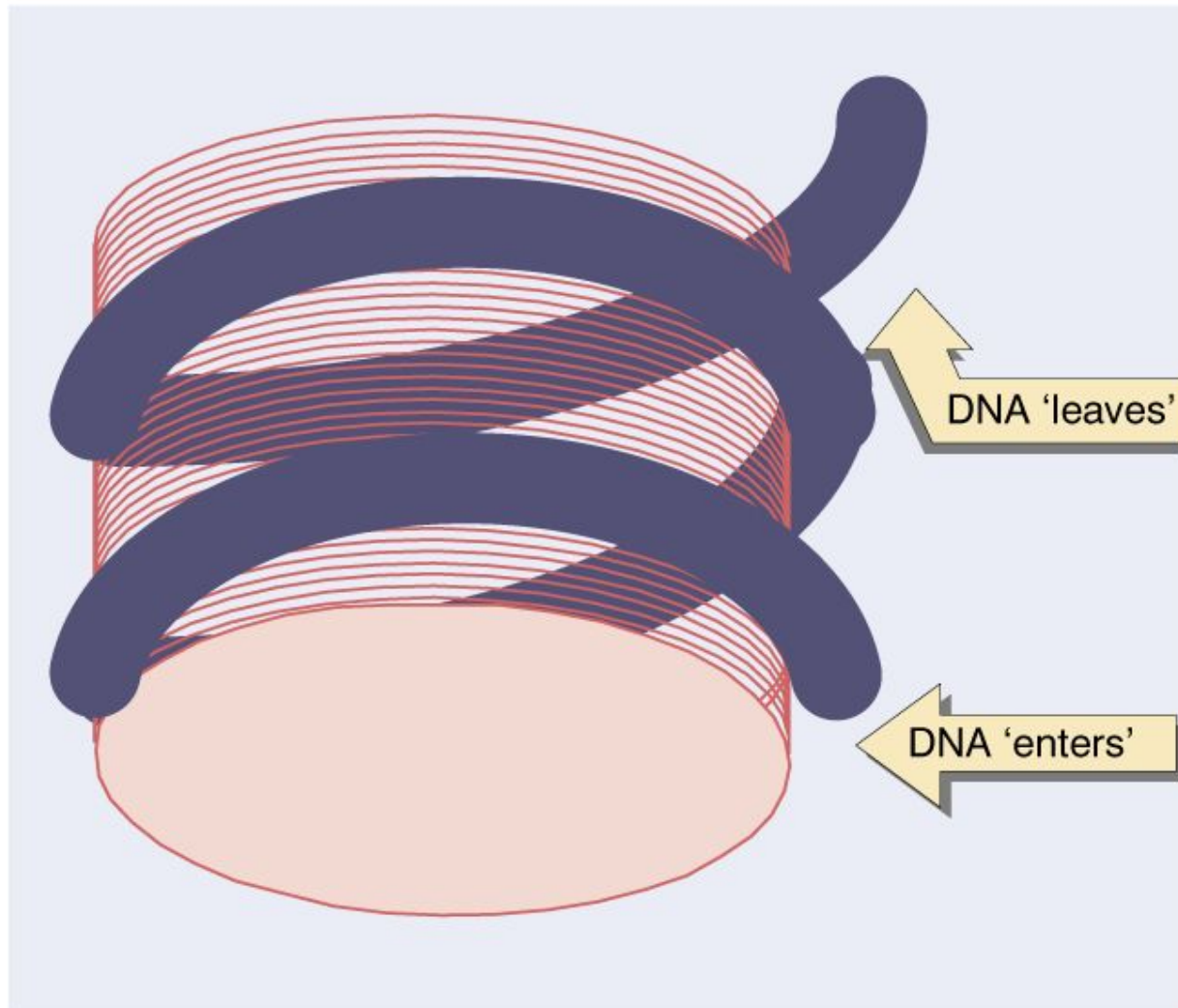
# Nucléosomes

**Figure 19.3** The nucleosome consists of approximately equal masses of DNA and histones (including H1). The predicted mass of the nucleosome is 262 kD.



# Nucléosomes

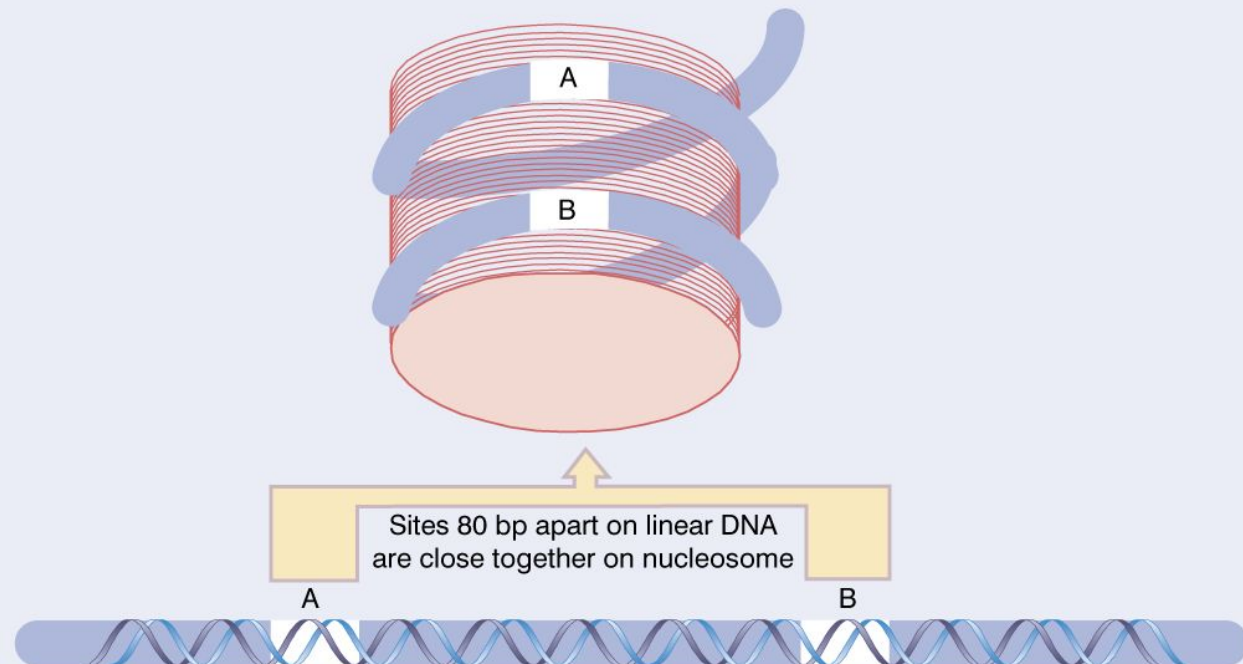
**Figure 19.4** The nucleosome may be a cylinder with DNA organized into two turns around the surface.



# Nucléosomes et organisation de l'ADN

**Figure 19.6**

Sequences on the DNA that lie on different turns around the nucleosome may be close together.

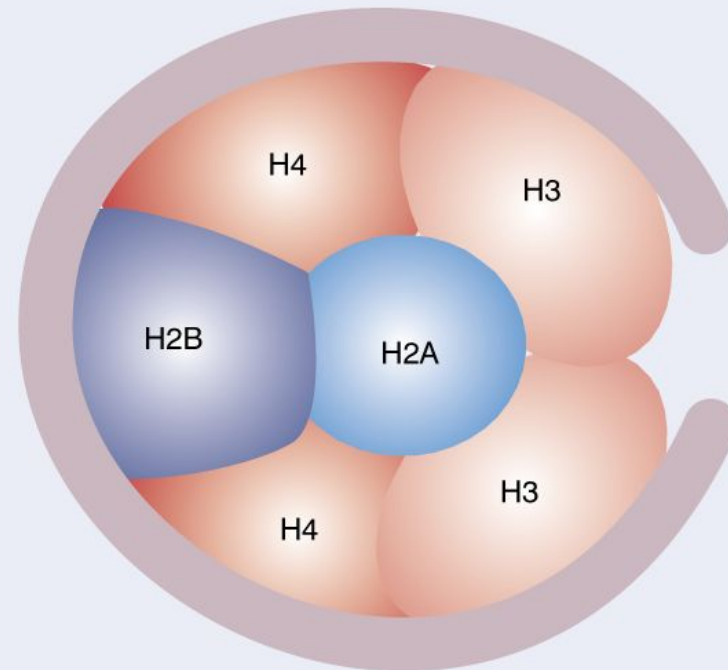


# Nucléosomes

**Figure 19.18** The 10 nm fiber is a continuous string of nucleosomes.



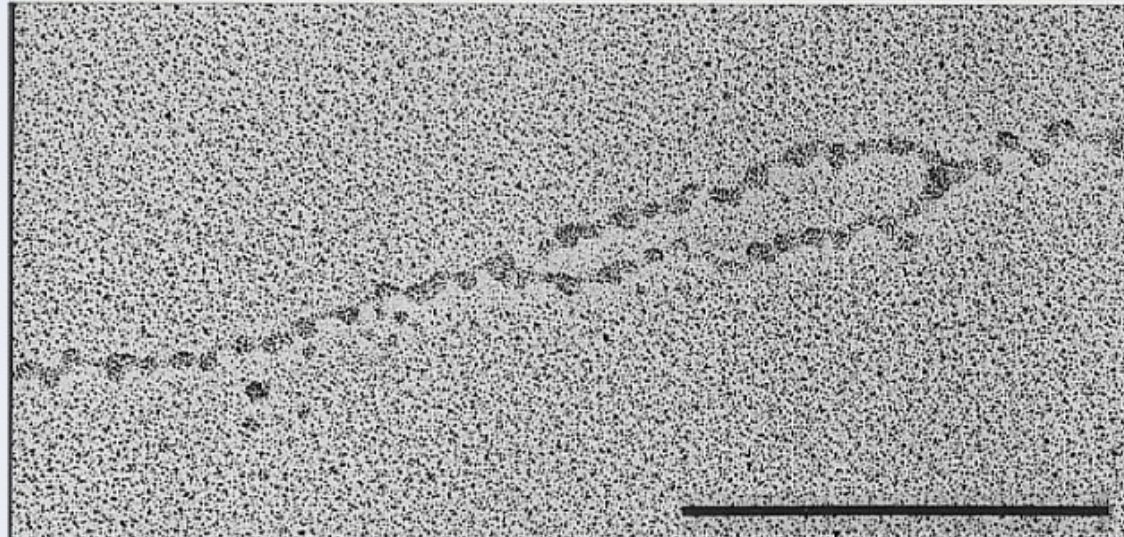
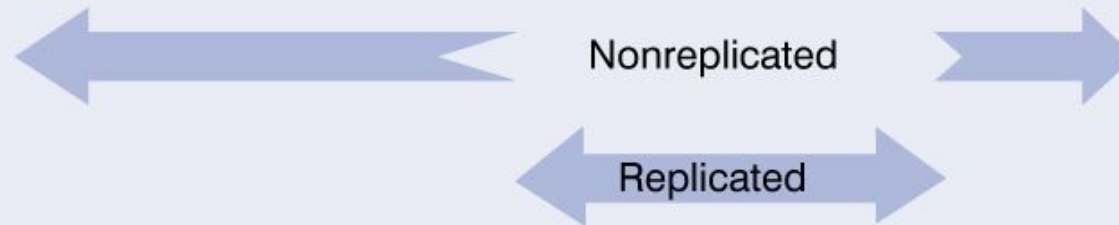
**Figure 19.21** In a symmetrical model for the nucleosome, the H3<sub>2</sub>-H4<sub>2</sub> tetramer provides a kernel for the shape. One H2A-H2B dimer can be seen in the top view; the other is underneath.





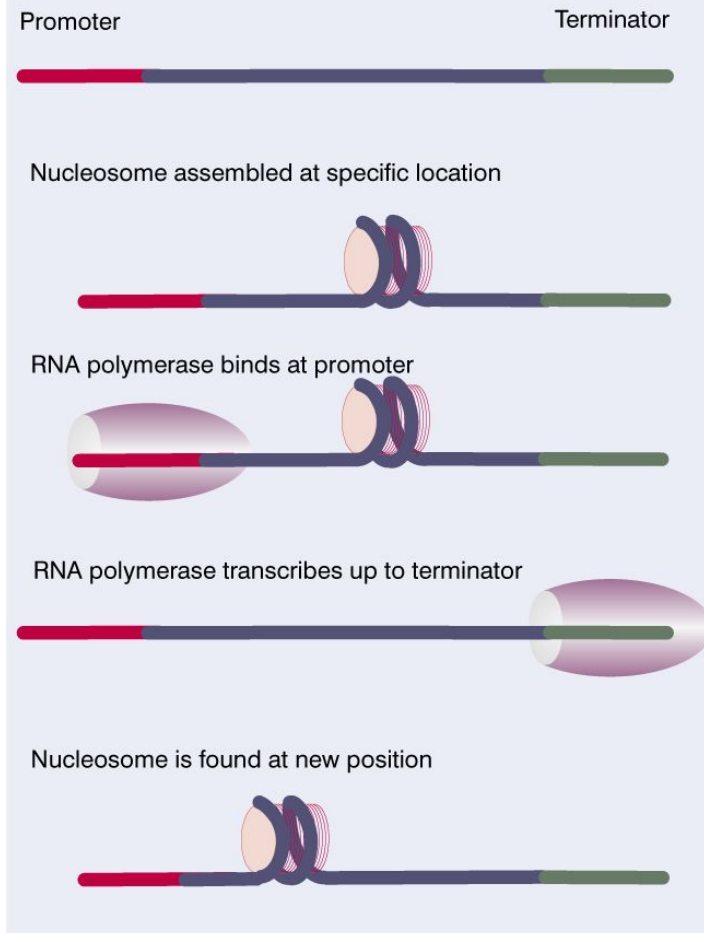
# Nucléosomes et réplication

**Figure 19.26** Replicated DNA is immediately incorporated into nucleosomes. Photograph kindly provided by S. MacKnight.



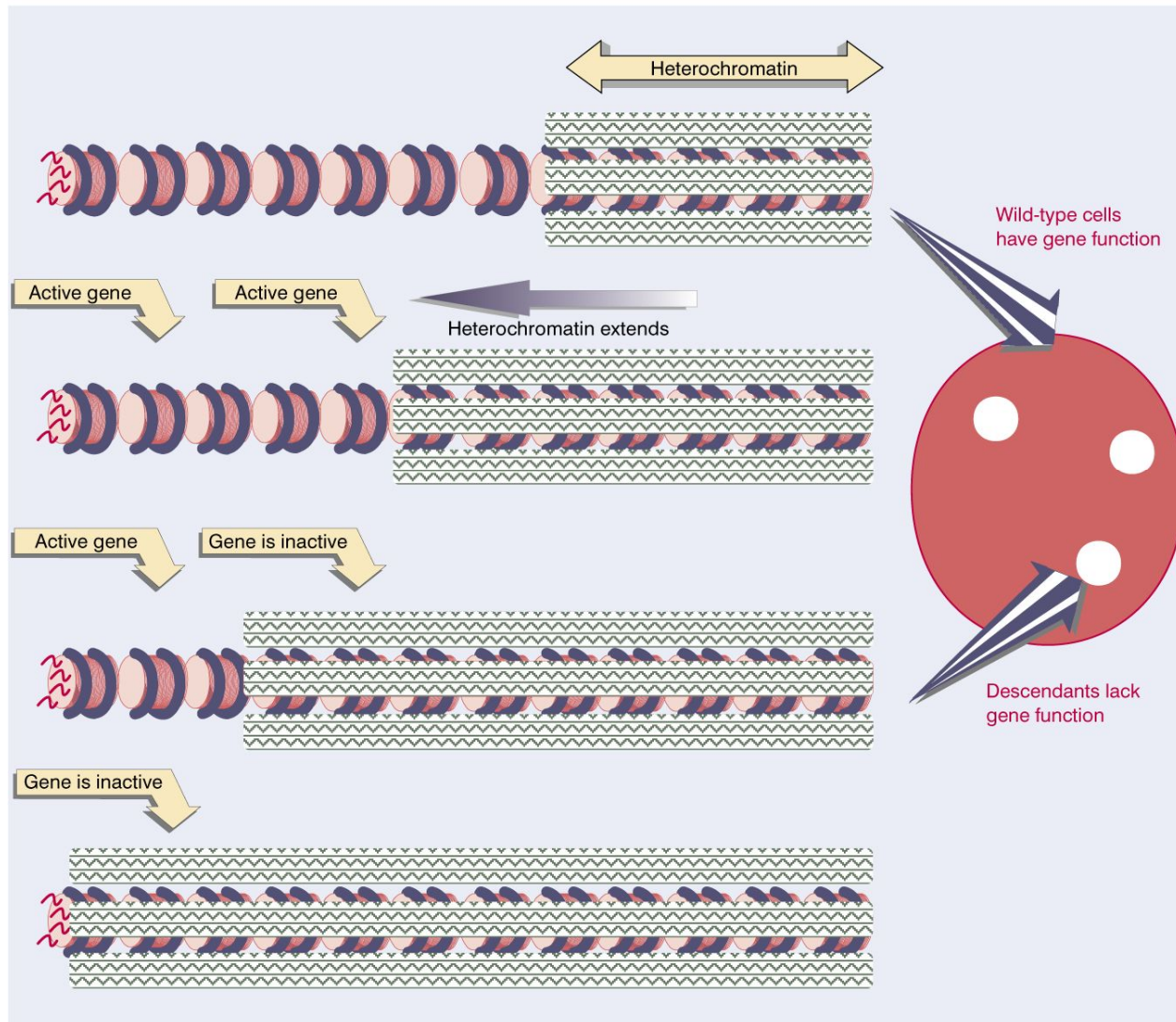
# Nucléosomes et transcription

**Figure 19.36** A protocol to test the effect of transcription on nucleosomes shows that the histone octamer is displaced from DNA and rebinds at a new position.



# Chromatine et inactivation génique

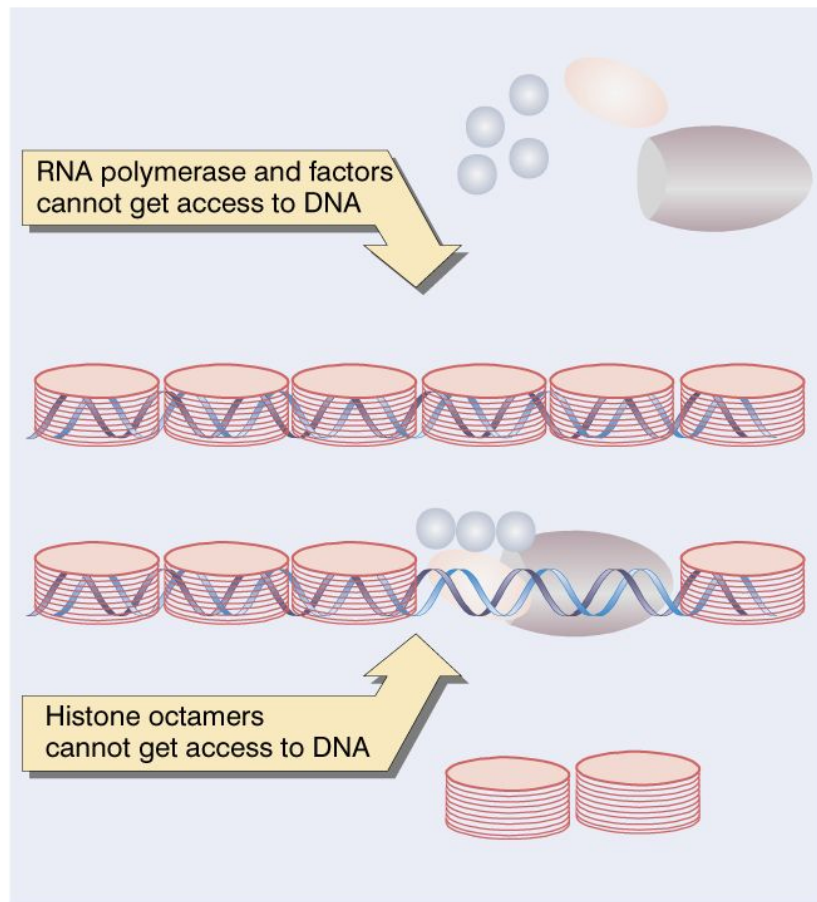
**Figure 19.45** Extension of heterochromatin inactivates genes. The probability that a gene will be inactivated depends on its distance from the heterochromatin region.



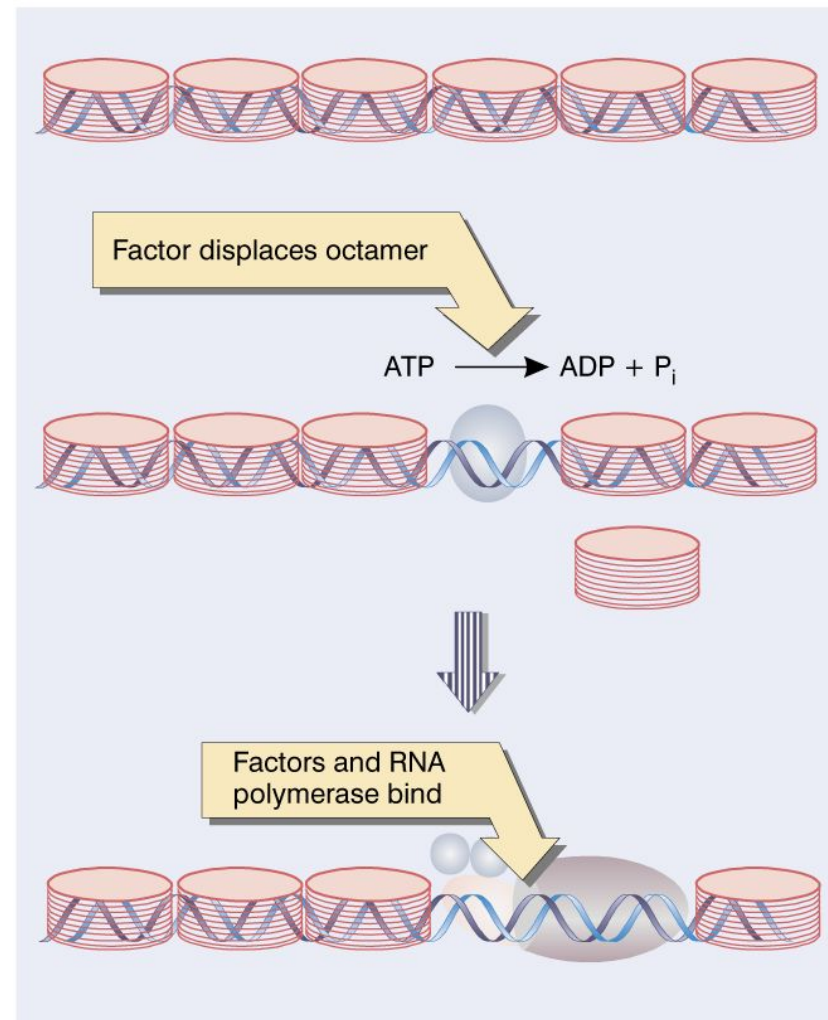


# Chromatine et transcription

**Figure 21.16** The pre-emptive model for transcription of chromatin proposes that if nucleosomes form at a promoter, transcription factors (and RNA polymerase) cannot bind. If transcription factors (and RNA polymerase) bind to the promoter to establish a stable complex for initiation, histones are excluded.



**Figure 21.17** The dynamic model for transcription of chromatin relies upon factors that can use energy provided by hydrolysis of ATP to displace nucleosomes from specific DNA sequences.



# Organisation génétique des génomomes

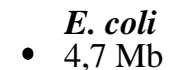
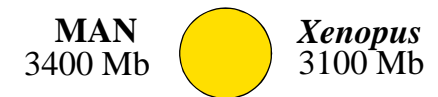
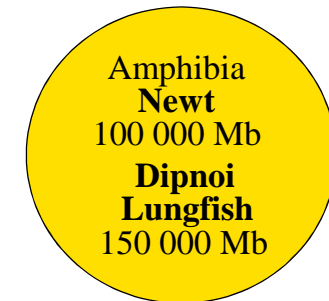
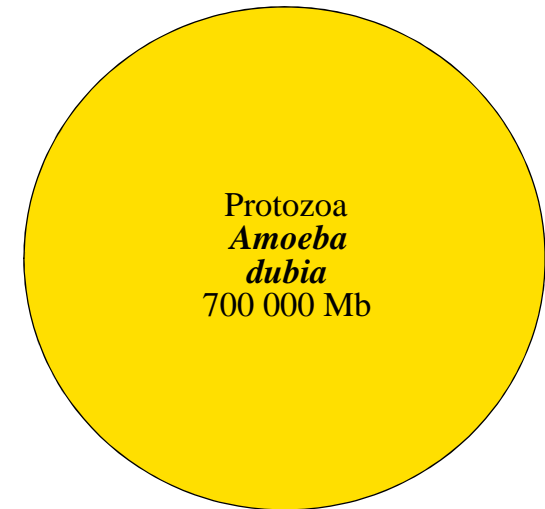
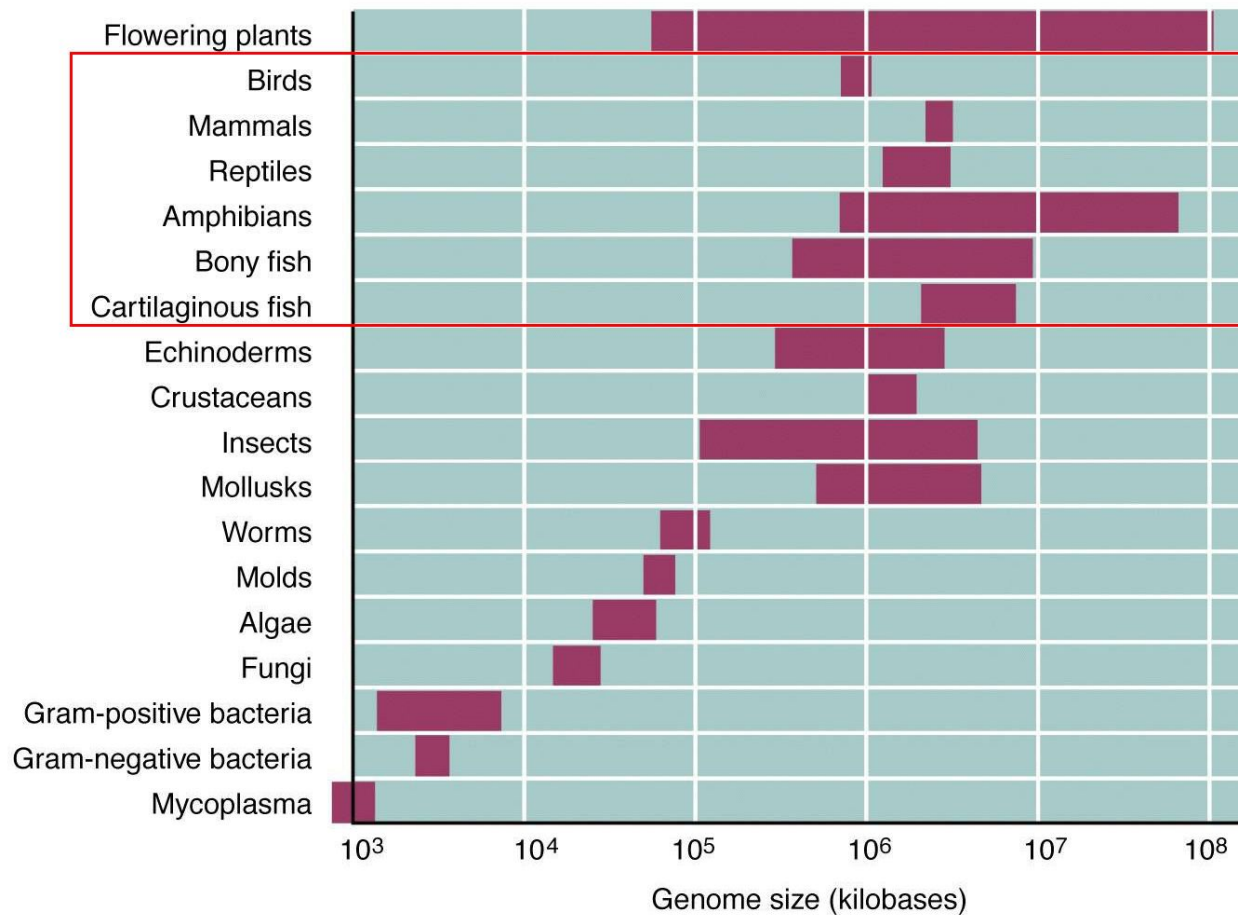
# Taille des génomes

- La longueur du génome et le nombre de gènes sont corrélés chez les procaryotes, pas chez les eucaryotes
- Le pourcentage de séquence codantes (densité) est très variable en fonction des organismes

**Figure 1.36** The amount of nucleic acid in the genome varies over an enormous range.

Genome	Gene Number	Base Pairs
<b>Organisms</b>		
Plants	<50,000	<10 <sup>11</sup>
Mammals	100,000	~3 × 10 <sup>8</sup>
Worms	14,000	~10 <sup>8</sup>
Flies	12,000	1.6 × 10 <sup>8</sup>
Fungi	6,000	1.3 × 10 <sup>7</sup>
Bacteria	2–4,000	<10 <sup>7</sup>
Mycoplasma	500	<10 <sup>6</sup>
<b>dsDNA Viruses</b>		
Vaccinia	<300	187,000
Papova (SV40)	~6	5,226
Phage T4	~200	165,000
<b>ssDNA Viruses</b>		
Parvovirus	5	5,000
Phage φX174	11	5,387
<b>dsRNA Viruses</b>		
Reovirus	22	23,000
<b>ssRNA Viruses</b>		
Coronavirus	7	20,000
Influenza	12	13,500
TMV	4	6,400
Phage MS2	4	3,569
STNV	1	1,300
<b>Viroids</b>		
PSTV RNA	0	359
<b>Scrapie</b>		
Prion	?	?

# Taille des génomes et complexité



# Gènes essentiels

- Un gène dont la délétion est létale est dit essentiel
- La recherche du sous ensemble minimal de gènes pour parvenir à faire vivre un organisme vivant est toujours d'actualité (environ 500 actuellement)

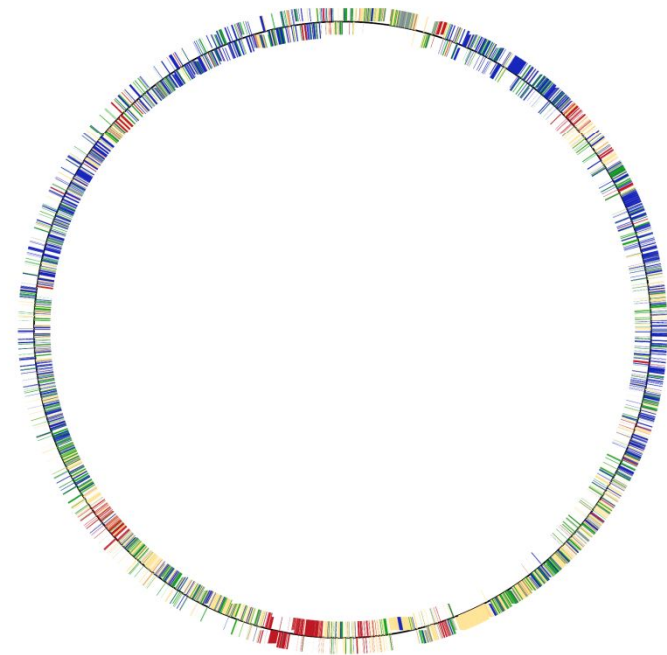
**Figure 3.10** Genome sizes, gene numbers and lethal loci.

Species	Genome (Mb)	Genes	Lethal loci
<i>Mycoplasma genitalium</i>	0.58	470	
<i>Rickettsia prowazekii</i>	1.11	834	
<i>Haemophilus influenzae</i>	1.83	1,743	
<i>Methanococcus jannaschi</i>	1.66	1,738	
<i>B. subtilis</i>	4.2	4,100	
<i>E. coli</i>	4.6	4,288	1,800
<i>S. cerevisiae</i>	13.5	6,034	3,600
<i>D. melanogaster</i>	165	12,000	3,100
<i>C. elegans</i>	97	19,099	
<i>H. sapiens</i>	3,300	100,000	

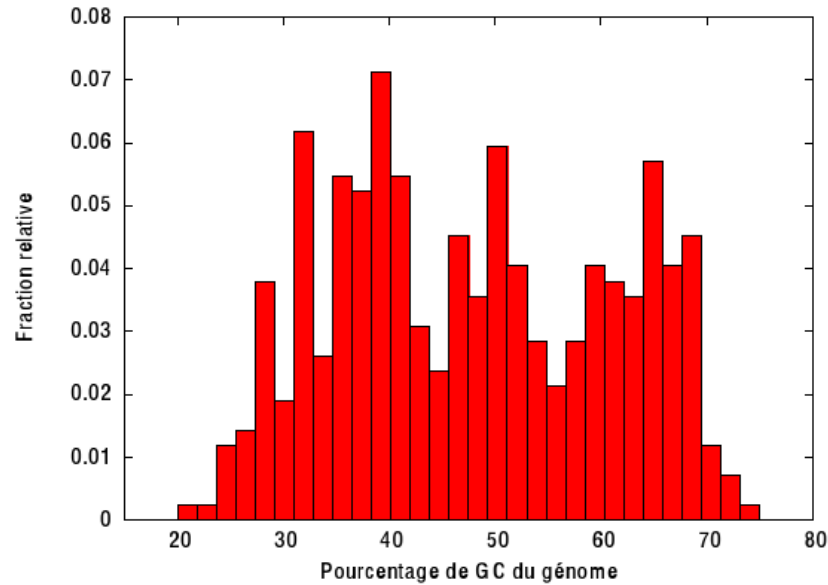


# Un génome bactérien : *E.coli*

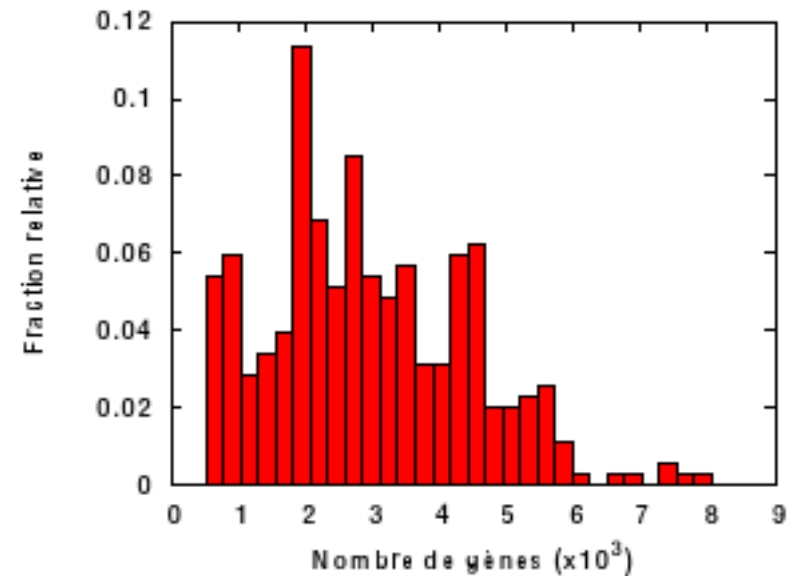
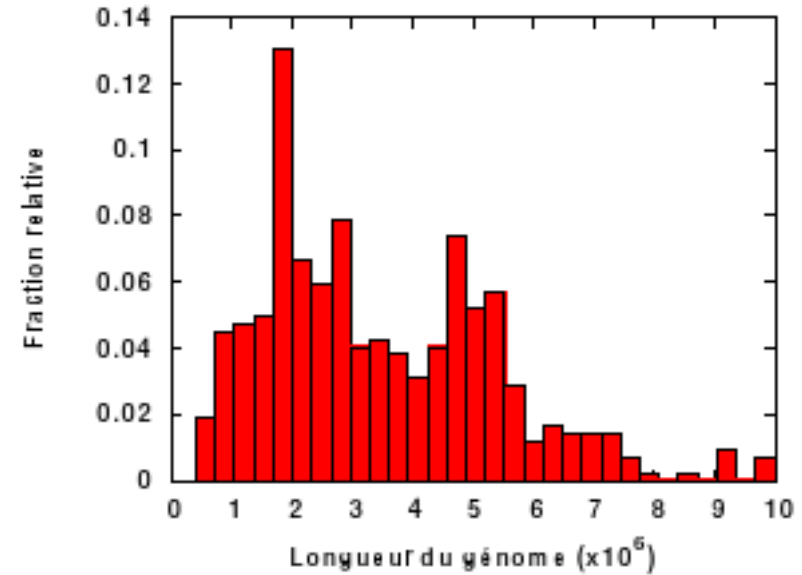
Genome Info:	Features:
Refseq: <a href="#">NC 000913</a>	Genes: <a href="#">4493</a>
GenBank: <a href="#">U00096</a>	Protein coding: <a href="#">4149</a>
Length: <b>4,639,675 nt</b>	Structural RNAs: <a href="#">172</a>
GC Content: <b>50%</b>	Pseudo genes: <b>177</b>
% Coding: <b>85%</b>	Others: <b>584</b>



# Diversité des génomes bactériens



Le nombre de plasmides des génomes bactériens varie de 0 à plus de 20 (dans ce cas 36% du génome)



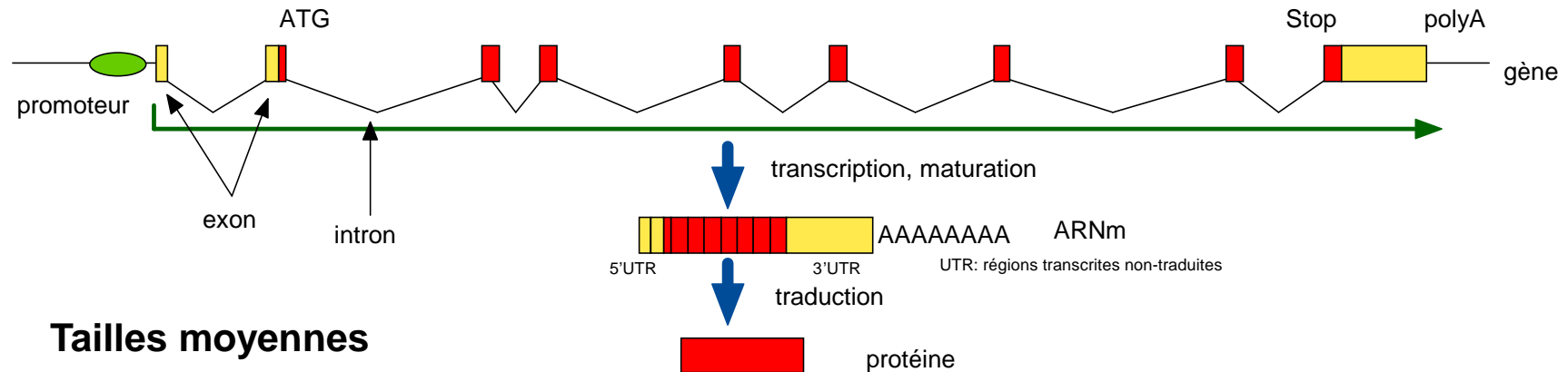
# Les contraintes en GC se traduisent au niveau protéique

- Les génomes riches ou pauvres en GC ont des possibilités restreintes en termes de séquences protéiques
- Le niveau de GC dans les génomes dépend de facteurs évolutifs et environnementaux

		Second letter				
		U	C	A	G	
First letter	U	UUU } Phe UUC } UUA } Leu UUG }	UCU } UCC } Ser UCA } UCG }	UAU } Tyr UAC } UAA Stop UAG Stop	UGU } Cys UGC } UGA Stop UGG Trp	U C A G
	C	CUU } CUC } Leu CUA } CUG }	CCU } CCC } Pro CCA } CCG }	CAU } His CAC } CAA } Gln CAG }	CGU } CGC } Arg CGA } CGG }	U C A G
	A	AUU } AUC } Ile AUA } AUG Met	ACU } ACC } Thr ACA } ACG }	AAU } Asn AAC } AAA } Lys AAG }	AGU } Ser AGC } AGA } Arg AGG }	U C A G
	G	GUU } GUC } Val GUA } GUG }	GCU } GCC } Ala GCA } GCG }	GAU } Asp GAC } GAA } Glu GAG }	GGU } GGC } Gly GGA } GGG }	U C A G



# Structure des gènes humains



- **Tailles moyennes**

– Gene	<b>45 kb</b>
– CDS	<b>1500 nt</b>
– Exon (interne)	<b>145 nt</b>
– Intron	<b>5200 nt</b>
– 5'UTR	<b>210 nt</b>
– 3'UTR	<b>740 nt</b>

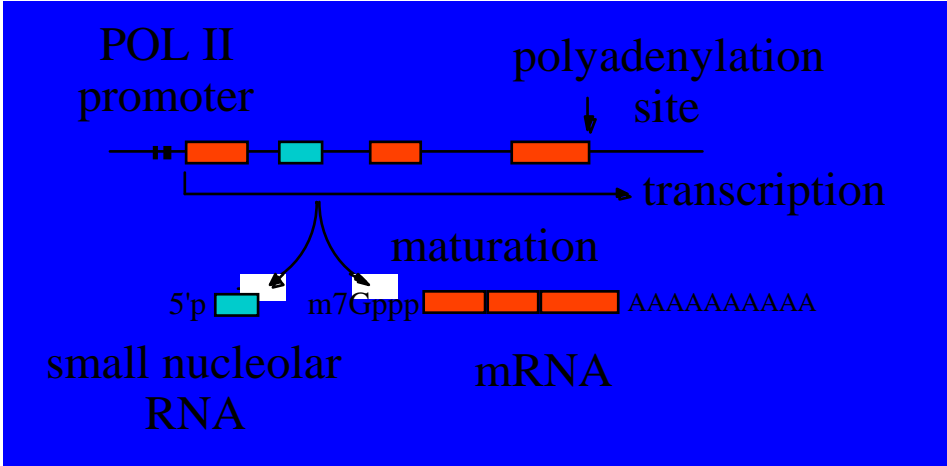
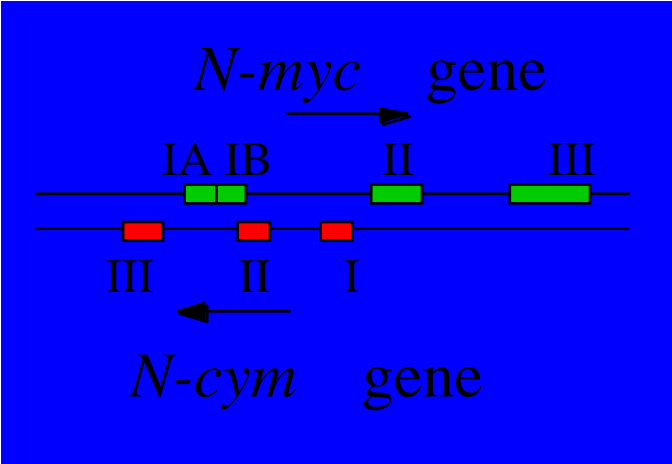
- **Intron / Exon**

– Nombres d'introns:	<b>6 ±3 introns / kb CDS</b>
– Introns / (introns + CDS):	<b>92%</b>

- **Epissage alternatif dans plus de 30% des gènes**

# Chevauchement?

Gènes partiellement superposés

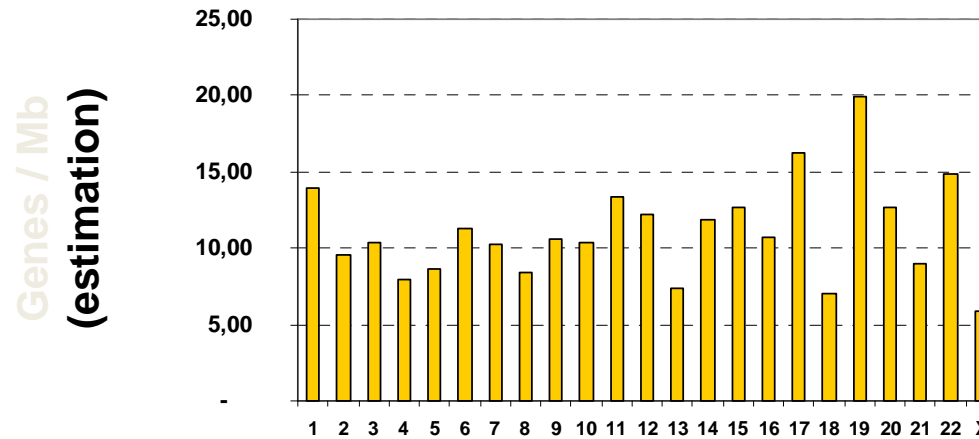
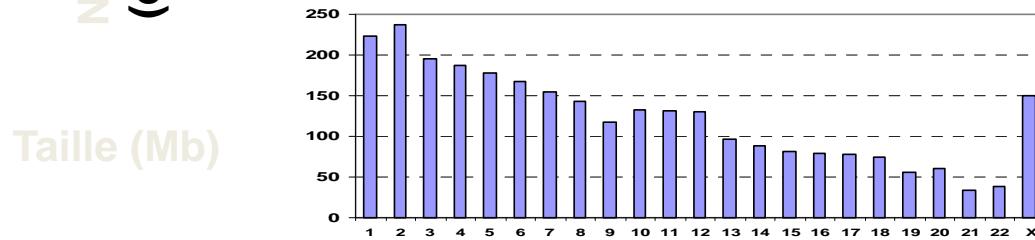
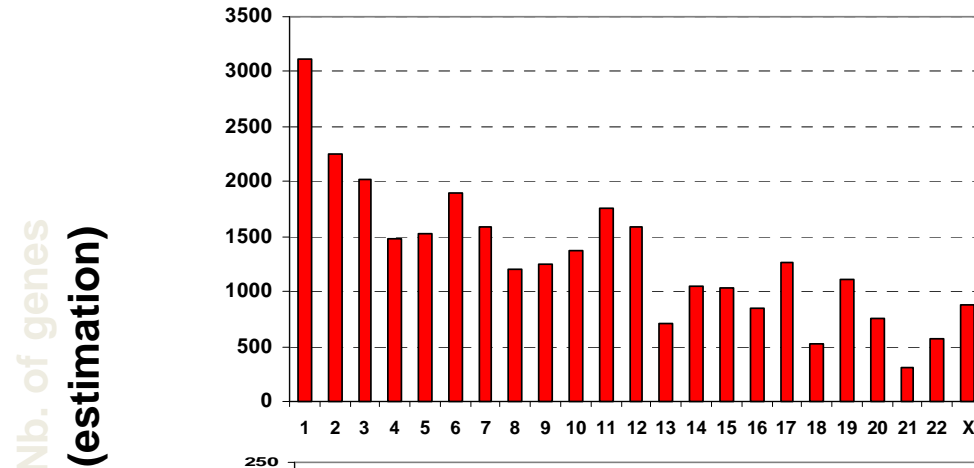


Petit ARN nucléolaire dans un intron de gène codant

# Gènes non traduits

- **tRNA (70-100 nt) :** ~ 350
- **rRNA :**
  - **18S (1.8kb), 5.8S (160 nt), 28S (5kb) :** 150-200
  - **5S (120 nt):** 200-300
- **snRNA (small nuclear RNA) (70-200 nt):** > 100  
Notamment épissage: U1, U2, U4, ...
- **snoRNA (small nucleolar RNA) (70-200 nt):** > 100  
Maturation des rRNA dans le nucléole
- **miRNA (micro RNA) :** 250 identifiés  
Régulation traduction, stabilité mRNA, transcription
- **ARN interférents:** ??

# Localisation des gènes



<http://www.ncbi.nlm.nih.gov/Science96/>

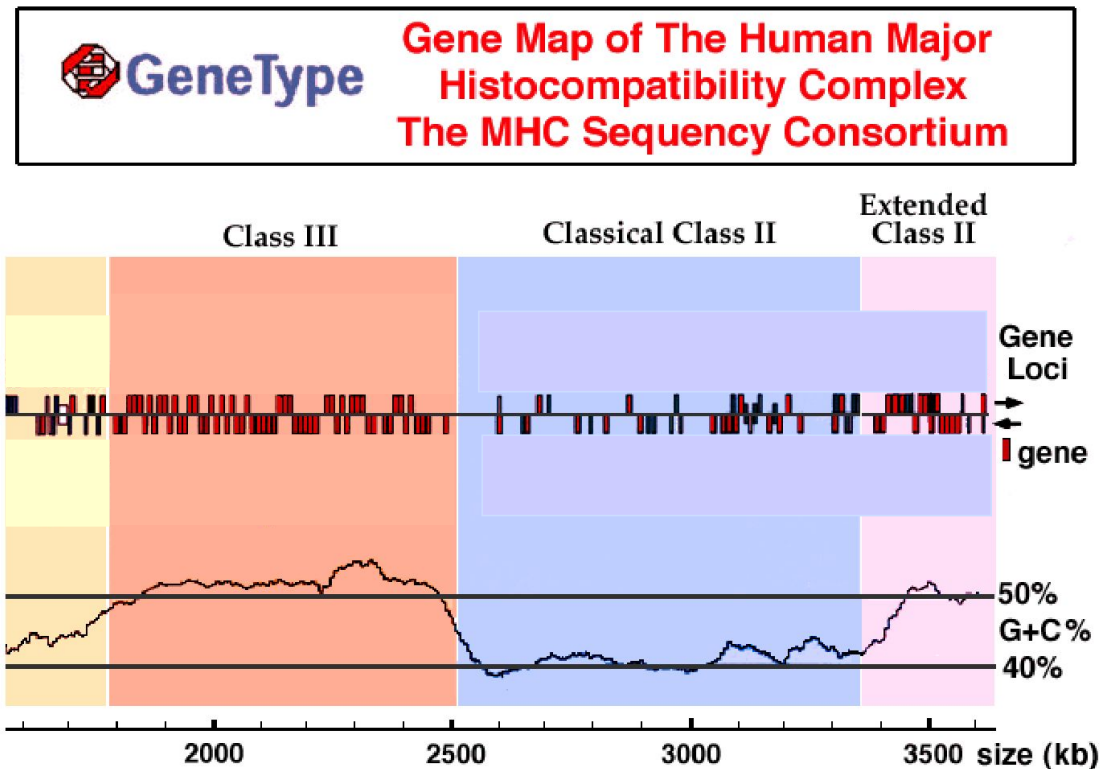


# Isochores

Génome de Vertébrés

Variations de la composition en bases le long des chromosomes

Sequence of human MHC



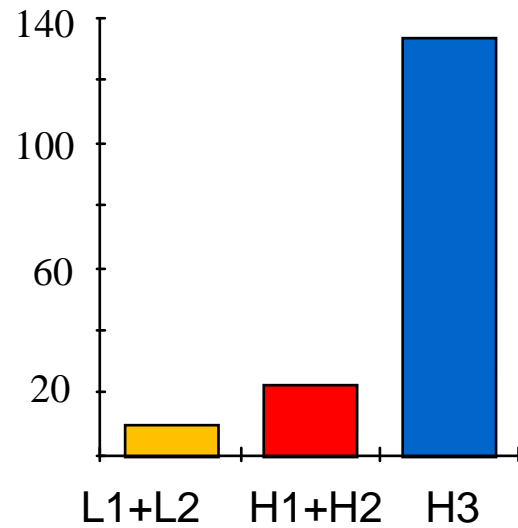
Définition

Régions > 300 kb homogène dans sa composition en bases

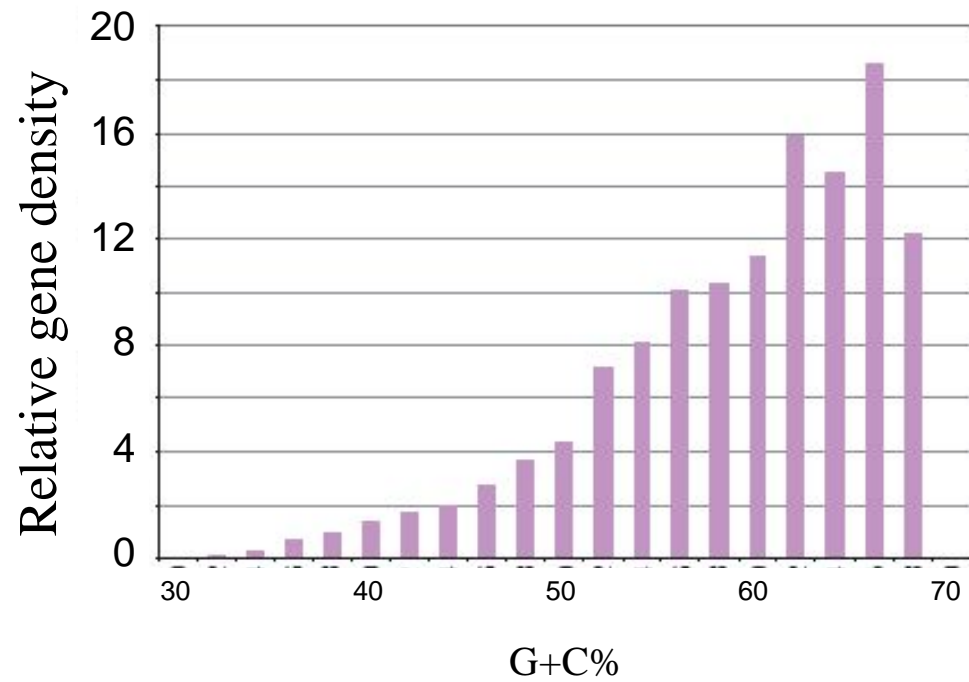
# Isochores

## Corrélation isochores riches en G+C / densité génique

Number of genes / Mb

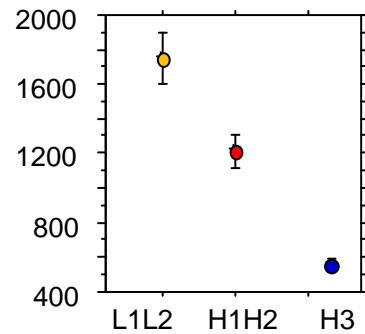


Mouchiroud et al. 1991

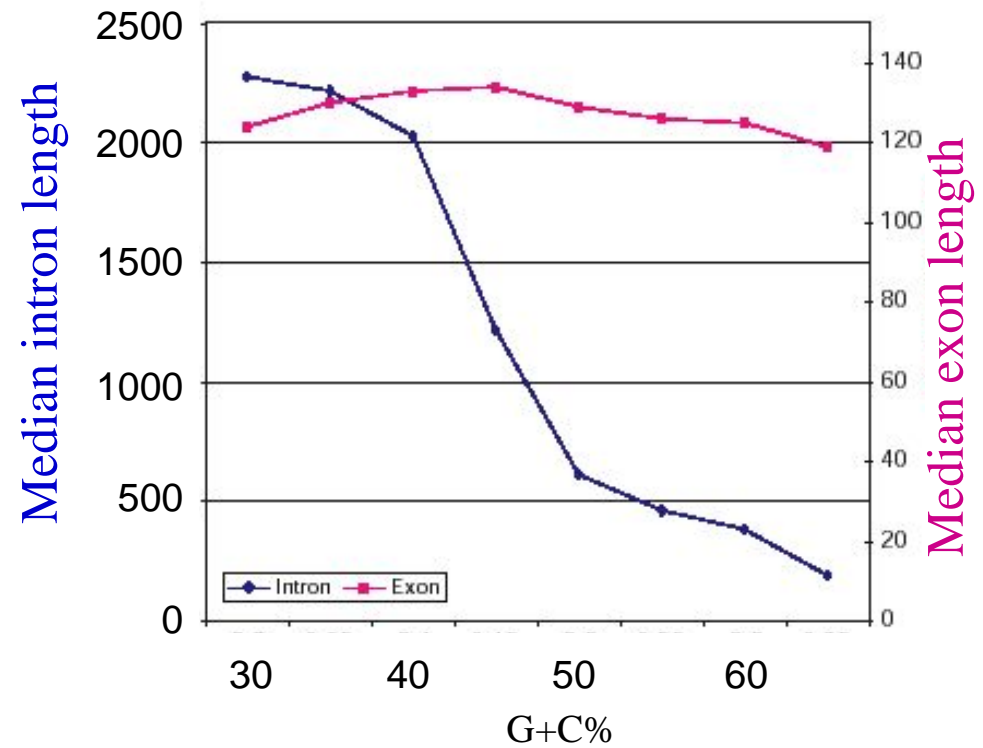
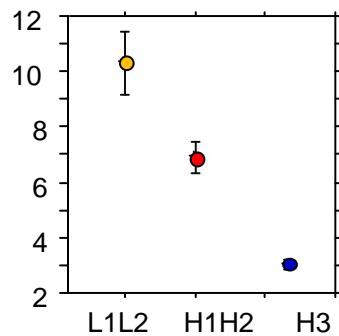


# Isochores et introns

Average intron length (bp)



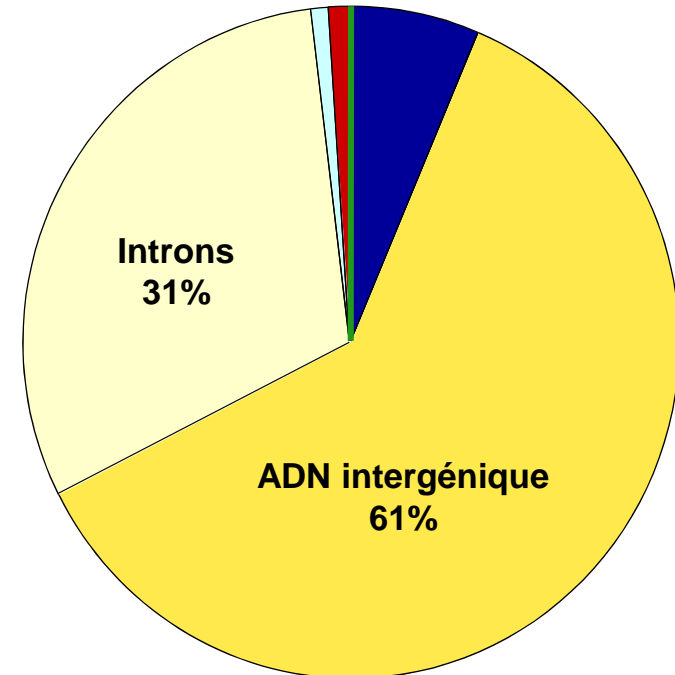
Gene compaction  
(intron length/coding region length)



→ Incidence sur la structure génétique

# Et l'ADN non codant?

- **Genes, regulatory elements**: ~ **2-5%**  
**20 000 – 25 000 proteic genes**  
**tRNA, rRNA, snRNA, miRNA**  
**UTR**
- **Non-coding sequences**: ~ **95-98%**
  - **Satellite DNA** (centromeres) ~ 6-7%
  - **Microsatellites** ~ 2%
  - **Transposable elements** ~ 46%
  - **Pseudogenes** ~ 1%
  - **Other (ancient -telements?)** ~ 42%
- **Variations in gene, repeat density and base composition (isochores) along chromosomes**



**>90% sans fonction connue**

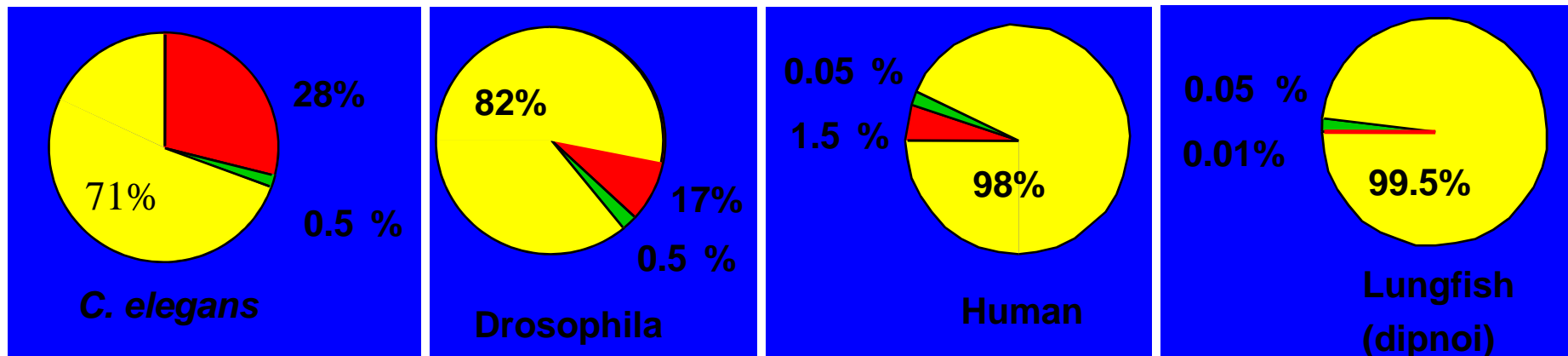
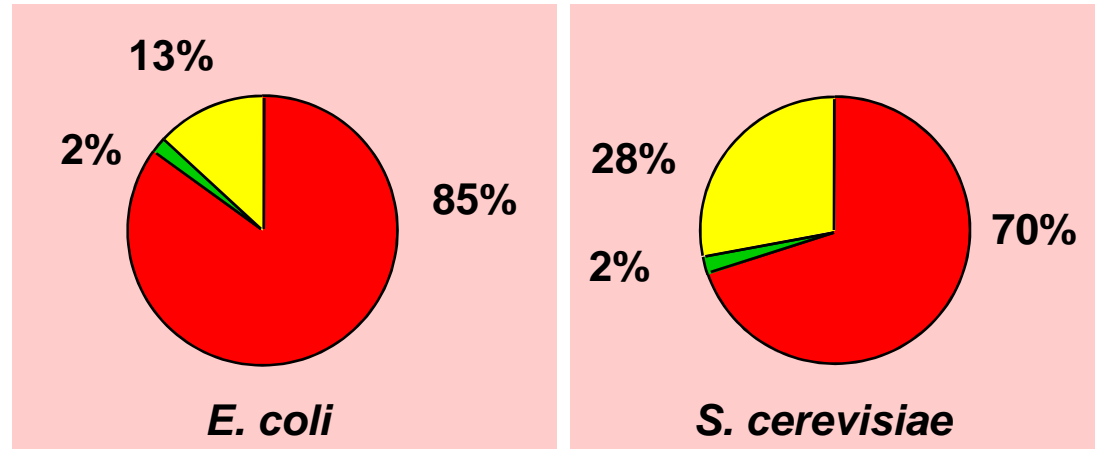
## **Elements fonctionnels non-quantifiés:**

- ARN non-traduits
- Promoteurs, enhancers
- ORI, MAR, télomères

## **Eléments non-fonctionnels:**

- Pseudogènes : 1.2%
- Elements transposables : 46%

# ADN non codant



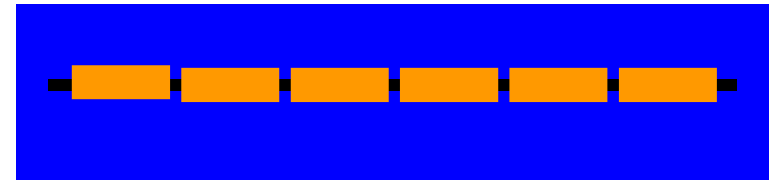
**Coding (protein)**      **RNA**      **Non-coding**

Proportion of functional elements within genomes

# Les séquences répétées

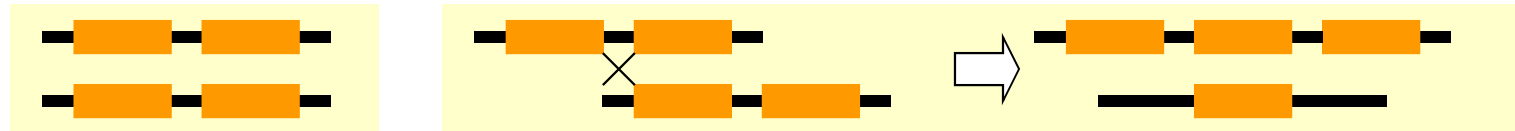
- Tandem repeats

	Motif(nt)	L(repeat)	%
Satellite	2-2000	→ 10Mb	6.5
Minisatellite	2-64	100-20000	0.3
Microsatellite	1-6	10-100	2



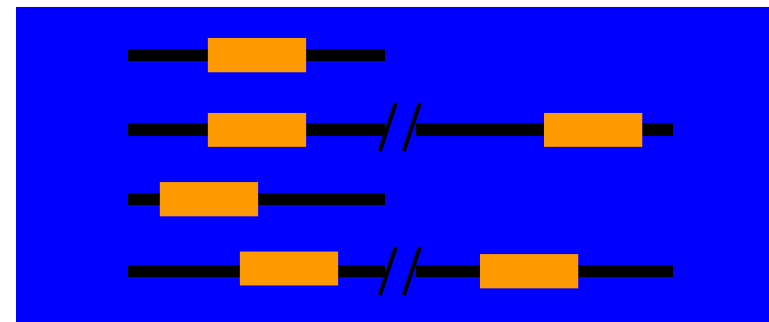
Slippage of the DNA-pol: CACACACACACA

Unequal crossing-over:



- Interspersed repeats

- DNA transposon (rares)
  - Retroelement
- } 46%



# Evolution des génomes

# Le modèle Jukes-Cantor

- Probabilité de substitution *par site par* unité de temps est  $\alpha$
- Substitution : 3 remplacements possibles  
(e.g.  $C \rightarrow \{A,G,T\}$ )
- Non-substitution: 1 seule possibilité  
(e.g.  $C \rightarrow C$ )

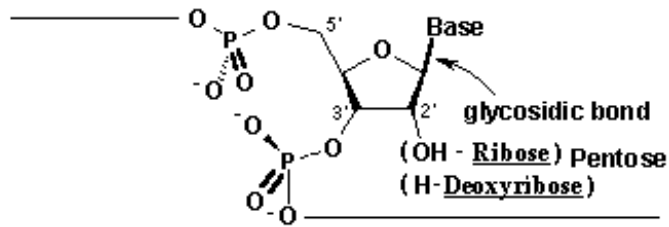


# Le modèle Jukes-Cantor

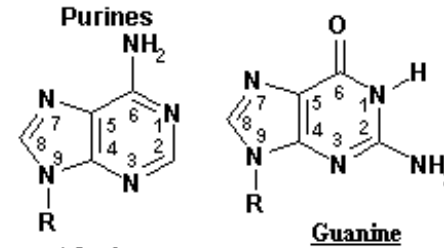
La matrice de transition de la chaîne de Markov a la forme suivante:

$$M_{JC} = \begin{array}{c|cccc} & \mathbf{A} & \mathbf{C} & \mathbf{G} & \mathbf{T} \\ \hline \mathbf{A} & 1-\alpha & \alpha/3 & \alpha/3 & \alpha/3 \\ \mathbf{C} & \alpha/3 & 1-\alpha & \alpha/3 & \alpha/3 \\ \mathbf{G} & \alpha/3 & \alpha/3 & 1-\alpha & \alpha/3 \\ \mathbf{T} & \alpha/3 & \alpha/3 & \alpha/3 & 1-\alpha \end{array}$$

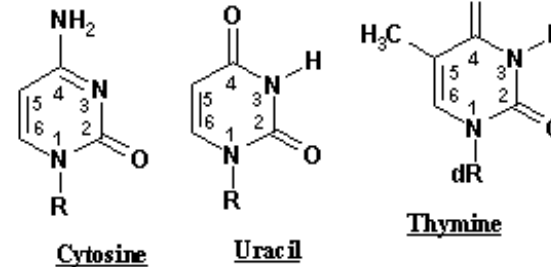
# Purines et pyrimidines



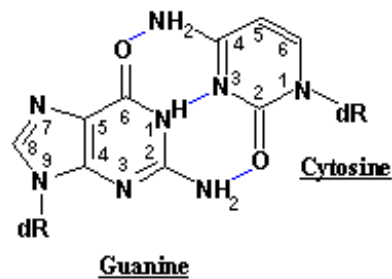
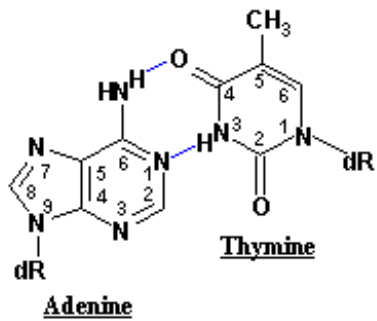
Nucleotide monophosphate



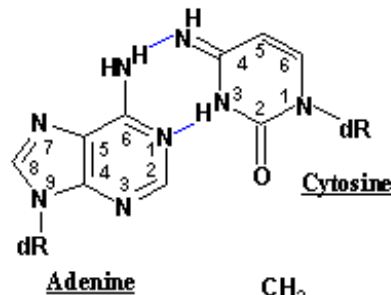
Pyrimidines



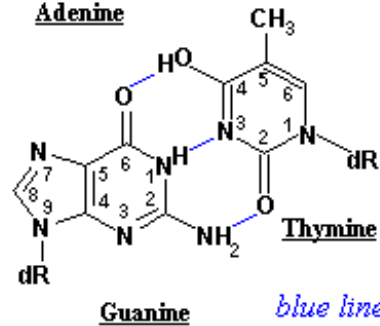
## CORRECT PAIRINGS



## WRONG PAIRINGS



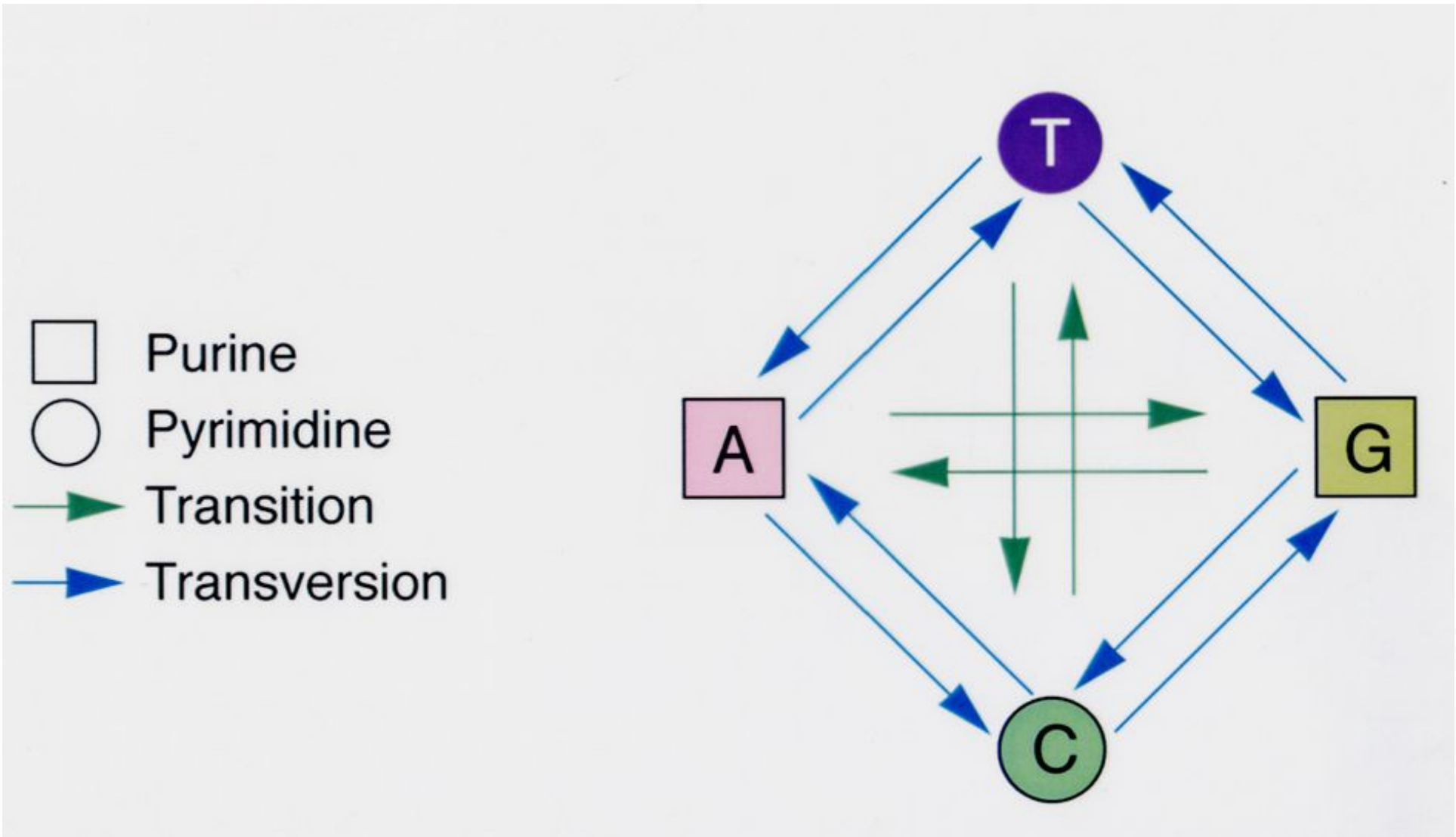
Cytosine in its imino form can pair wrongly with adenine



Thymine in its enol form can pair wrongly with guanine

blue lines = hydrogen bonds

# Transitions et transversions



# Le modèle de Kimura

Inclut le biais de substitution entre transitions et transversions

**Probabilité de transition** ( $G \leftrightarrow A$  et  $T \leftrightarrow C$ ) *par site par* unité de temps est  $\alpha$

**Probabilité de transversion** ( $G \leftrightarrow T$ ,  $G \leftrightarrow C$ ,  $A \leftrightarrow T$ , and  $A \leftrightarrow C$ ) *par site par* unité de temps est  $\beta$

# Le modèle de Kimura

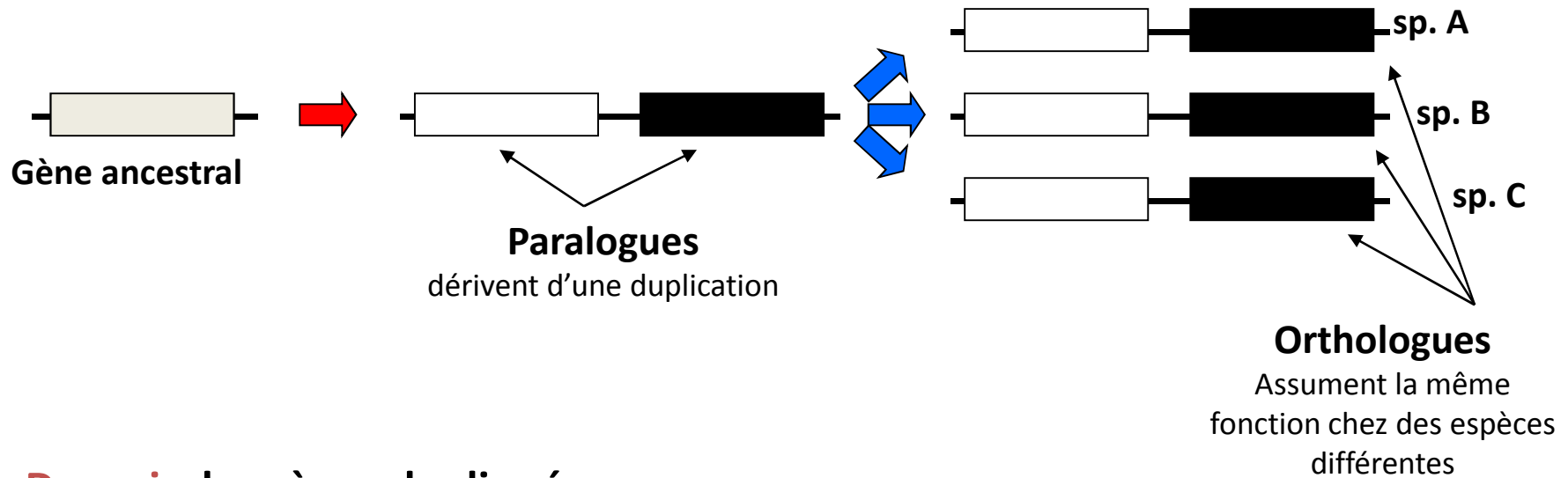
La matrice de transition de la chaîne de Markov a alors la forme suivante:

$$M_{K2P} = \begin{array}{ccccc} & \mathbf{A} & \mathbf{C} & \mathbf{G} & \mathbf{T} \\ \mathbf{A} & 1-\alpha-\beta & \beta & \alpha & \beta \\ \mathbf{C} & \beta & 1-\alpha-\beta & \beta & \alpha \\ \mathbf{G} & \alpha & \beta & 1-\alpha-\beta & \beta \\ \mathbf{T} & \beta & \alpha & \beta & 1-\alpha-\beta \end{array}$$

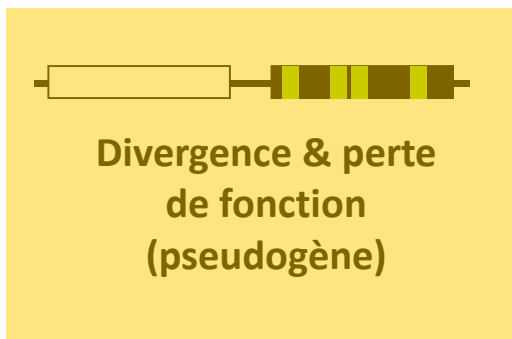
# Modèles de substitution plus avancés

- Plus réalistes
- Le modèle de Kimura a une distribution uniforme stationnaire comme solution, ce qui n'est pas vrai
- Une possibilité: la probabilité de substitution d'une purine à une pyrimidine est différente de celle d'une pyrimidine à une purine (solution dynamique).
- Le modèle HKY capture ce biais

# Duplications géniques



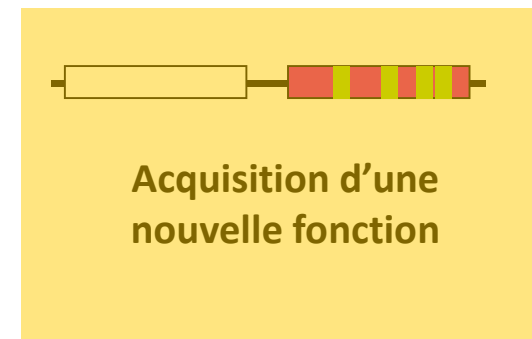
## Devenir des gènes dupliqués



*PSEUDOGENISATION*



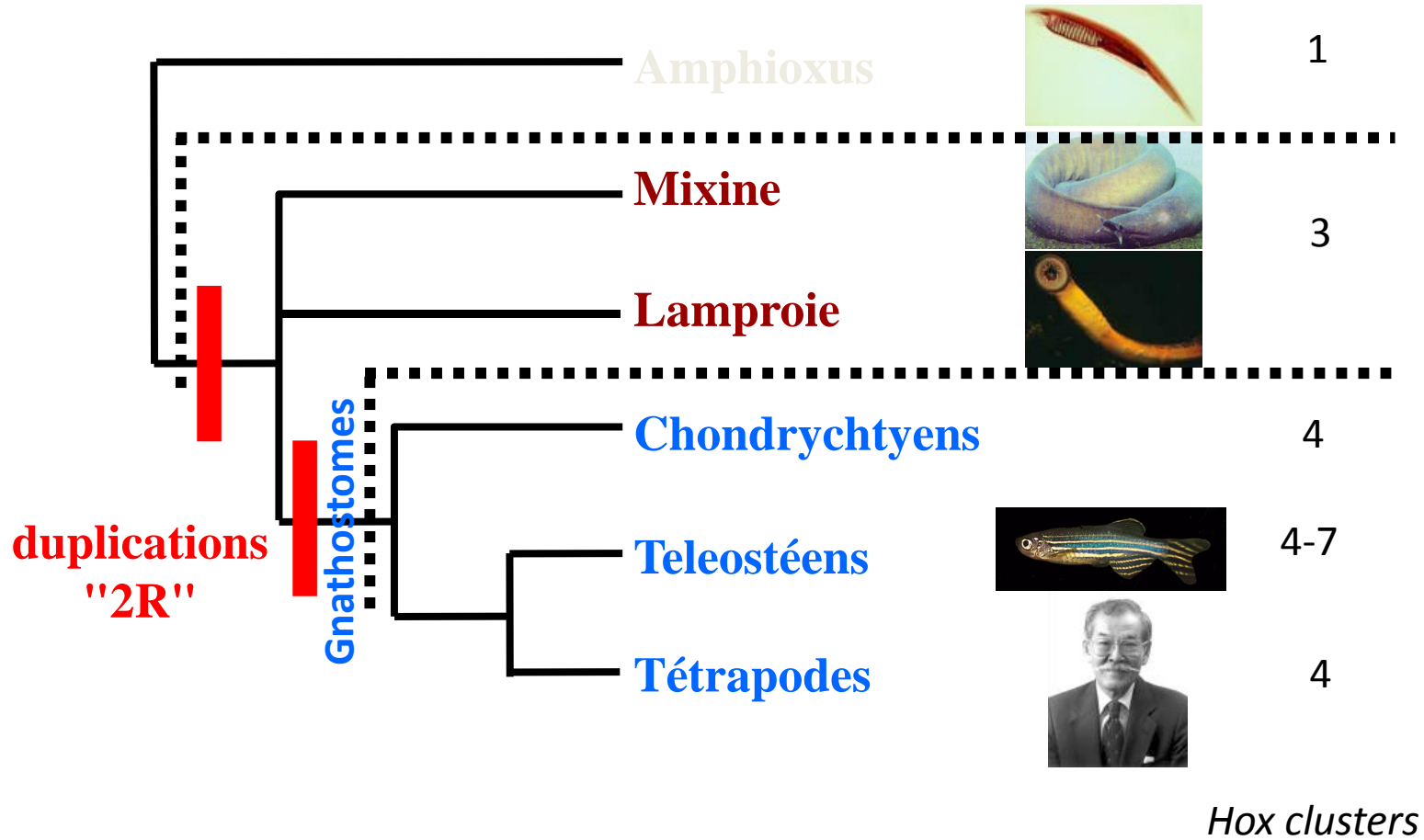
*SUB-FONCTIONNALISATION*



*NEO-FONCTIONNALISATION*

# Duplications de groupes de gènes

Evolution des espèces





# Duplications, paralogie et phylogénie

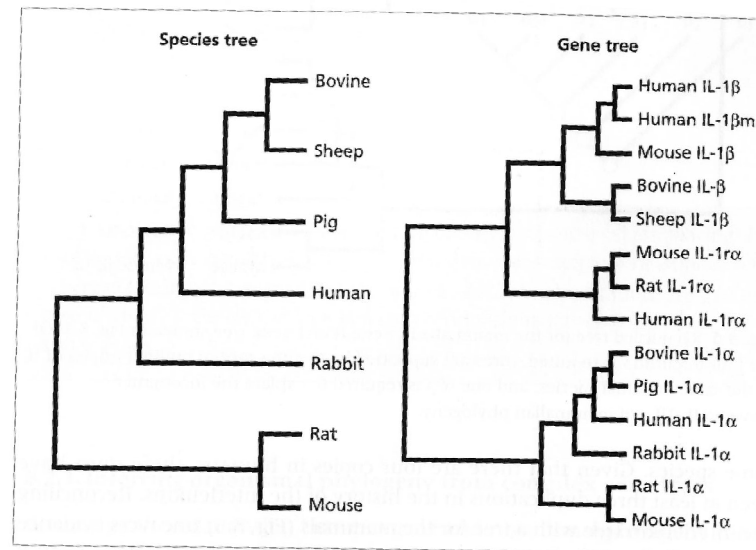
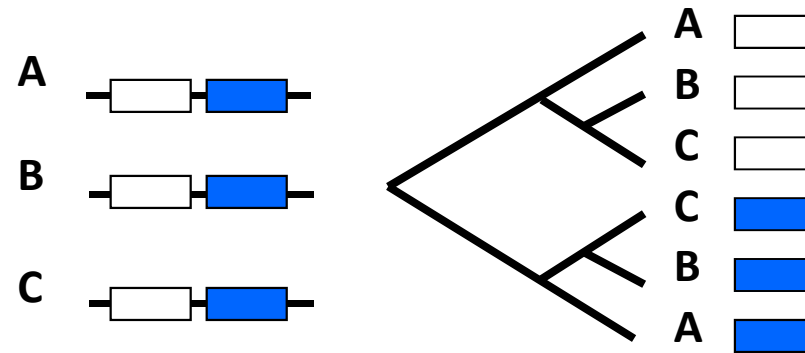


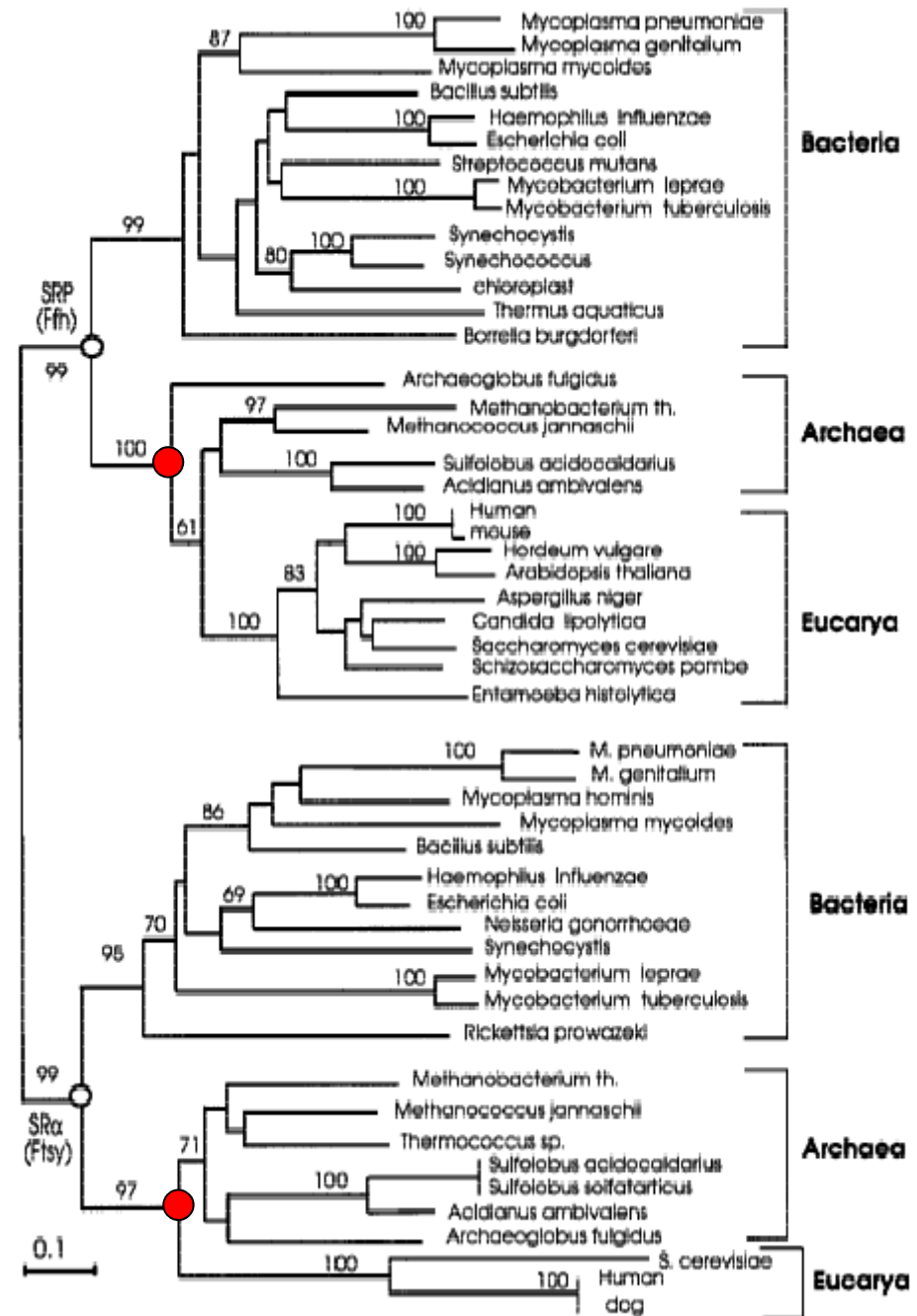
Fig. 8.5 Tree for mammals and for mammalian interleukin-1 genes.

# The Root of the Universal Tree of Life Inferred from Anciently Duplicated Genes Encoding Components of the Protein-Targeting Machinery

Simonetta Gribaldo, Piero Cammarano

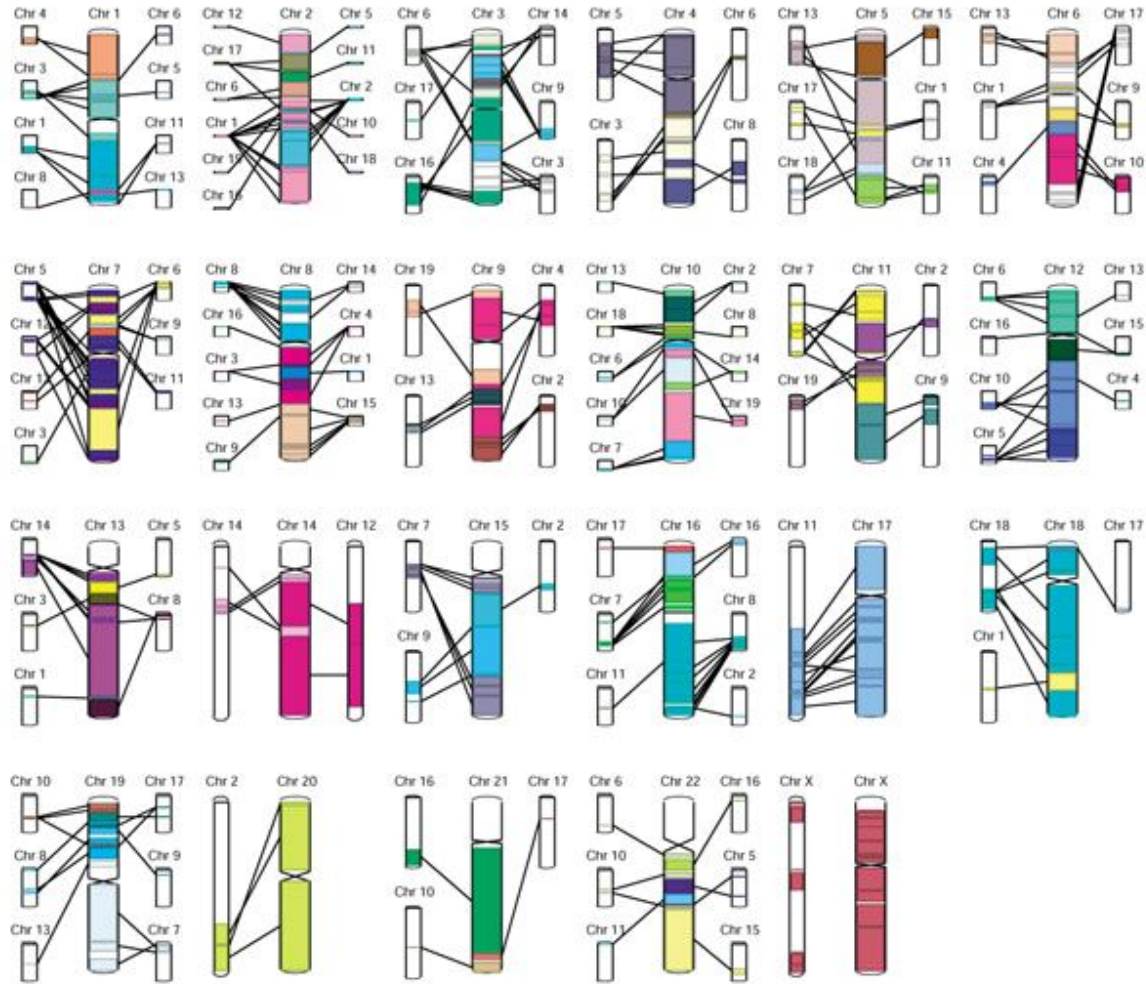
J Mol Evol (1998) 47:508–516

**Fig. 3.** Neighbor-joining tree constructed from the SRP54(Ffh)/SR $\alpha$ (Ftsy) alignment. Numbers attached to internal nodes are BCL based on 100 bootstrap replicates of the original alignment. BCLs <50% are not shown. Evolutionary distance matrices were generated by the PROTDIST program with the "Dayhoff" option (see Methods). Scale bar represents 0.1 amino acid substitution per site.



ACQUISITION DU NOYAU

# Synténie

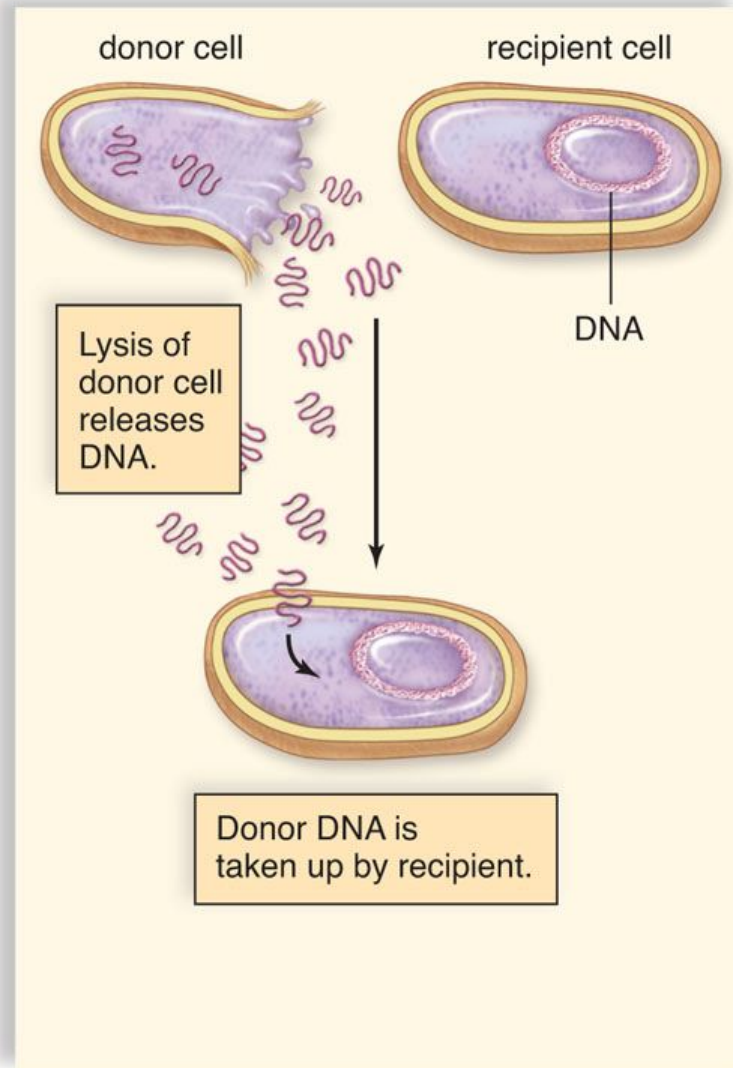


# Transfert horizontal

- **Transformation** – une bactérie ingère de l'ADN provenant du milieu environnant.
- **Conjugaison** – La bactérie donneuse envoie de l'ADN à une autre bactérie par l'intermédiaire d'un pilus sexuel.
- **Transduction** – Un bactériophage transfère une portion d'ADN d'une bactérie à l'autre au cours d'infections multiples.

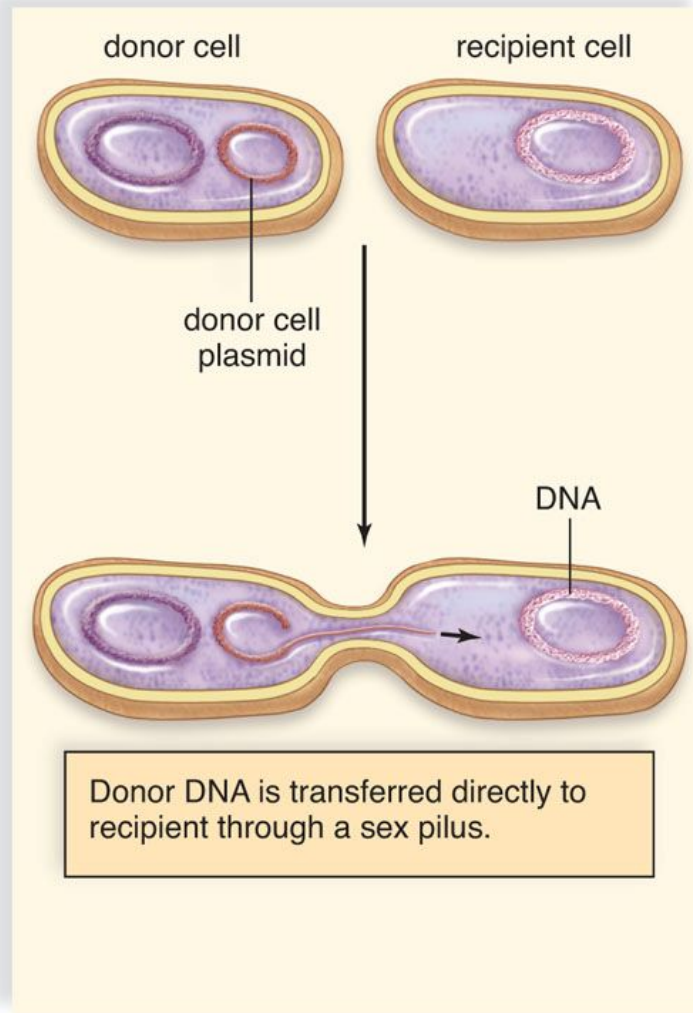
# Transformation

Copyright © The McGraw-Hill Companies, Inc. Permission required for reproduction or display.



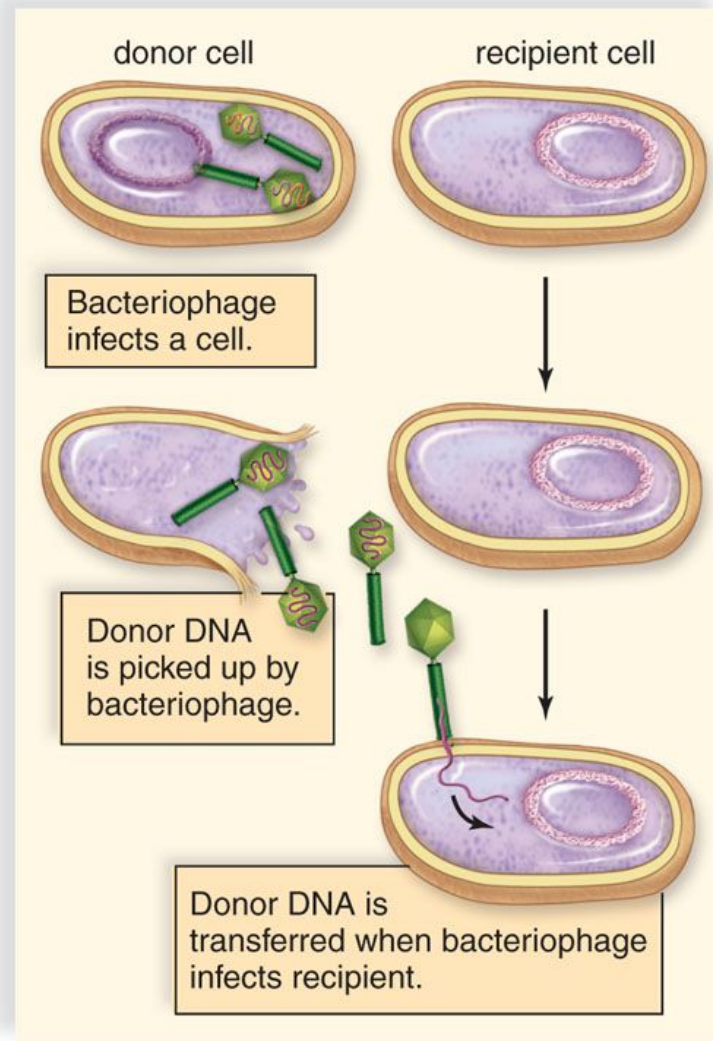
# Conjugaison

Copyright © The McGraw-Hill Companies, Inc. Permission required for reproduction or display.



# Transduction

Copyright © The McGraw-Hill Companies, Inc. Permission required for reproduction or display.

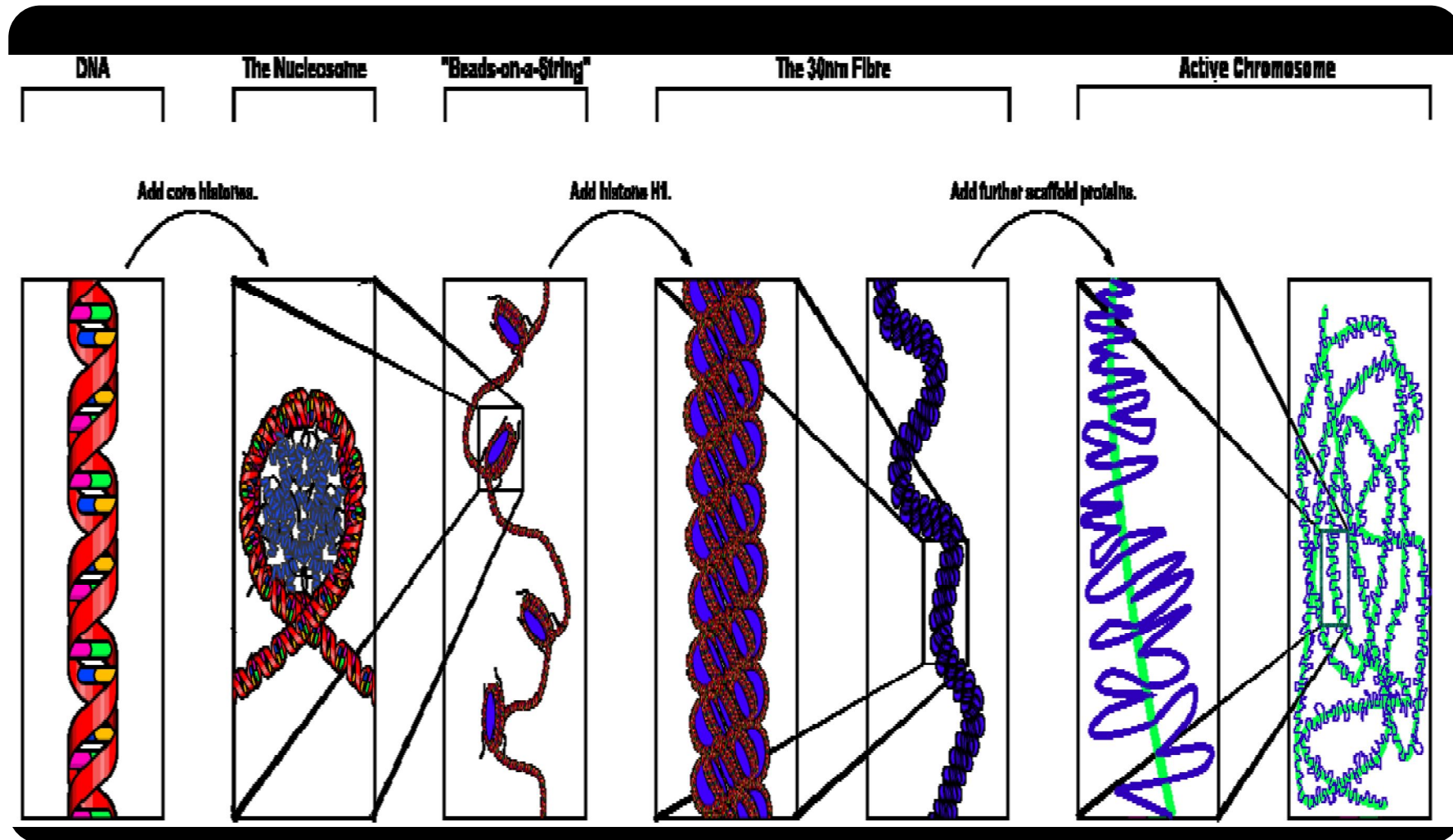








# Chromosome Packing



**Figure 3.15** The chloroplast genome codes for 4 rRNAs, 30 tRNAs, and ~50 proteins.

Genes	Types
<b>RNA-coding</b>	
16S rRNA	1
23S rRNA	1
4.5S rRNA	1
5S rRNA	1
tRNA	30
<b>Gene Expression</b>	
r-proteins	19
RNA polymerase	3
Others	2
<b>Thylakoid Membranes</b>	
Photosystem I	2
Photosystem II	7
Cytochrome <i>b/f</i>	3
H <sup>+</sup> -ATPase	6
<b>Others</b>	
NADH dehydrogenase	6
Ferredoxin	3
Ribulose BP Cblase	1
Unidentified	29
<b>Total</b>	<b>110</b>

# Root of the universal tree of life based on ancient aminoacyl-tRNA synthetase gene duplications

(Archaea/Bacteria/eukaryote/phylogeny)

JAMES R. BROWN\* AND W. FORD DOOLITTLE

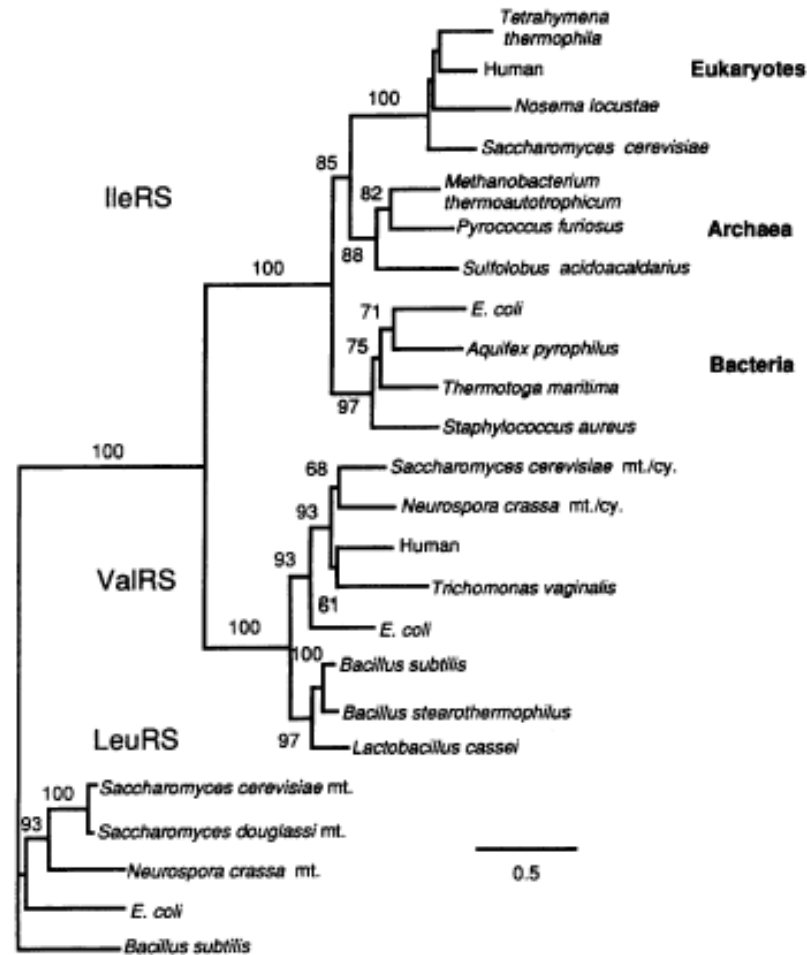


FIG. 3. Neighbor-joining tree of IleRS, ValRS, and LeuRS genes using the program NEIGHBOR (26). The scale represents 0.5 expected number of amino acid replacements per position as determined with the program PROTDIST. Numbers are the frequency of occurrence of nodes that exceeded 50% of 300 bootstrap replicates.

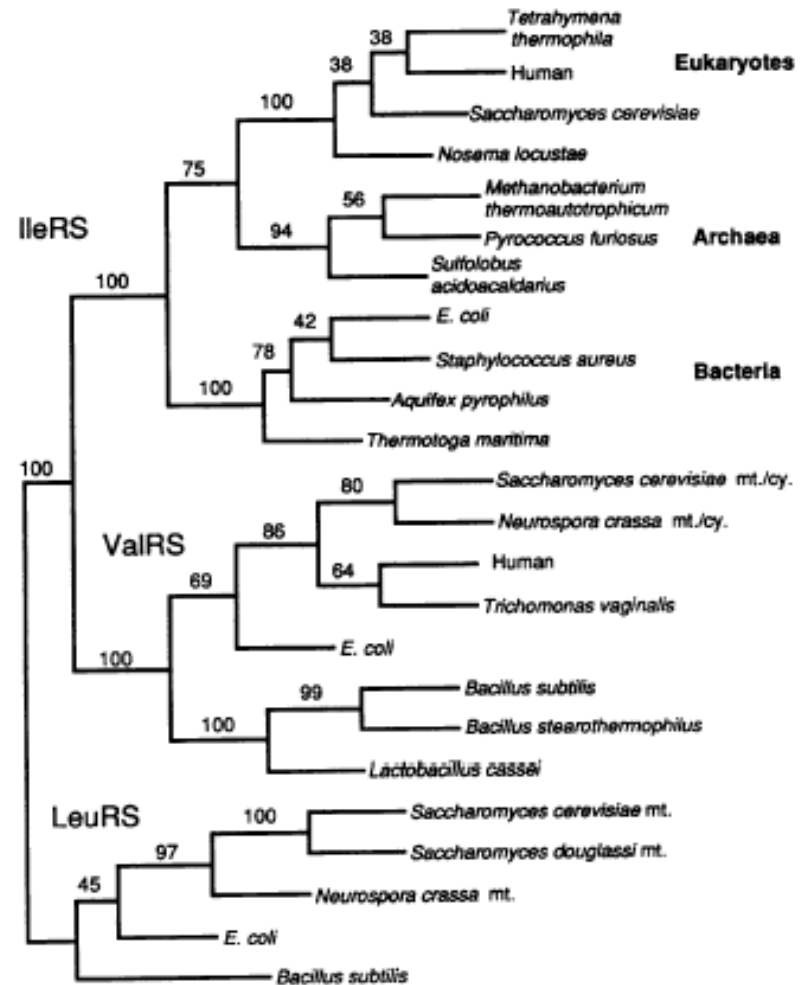
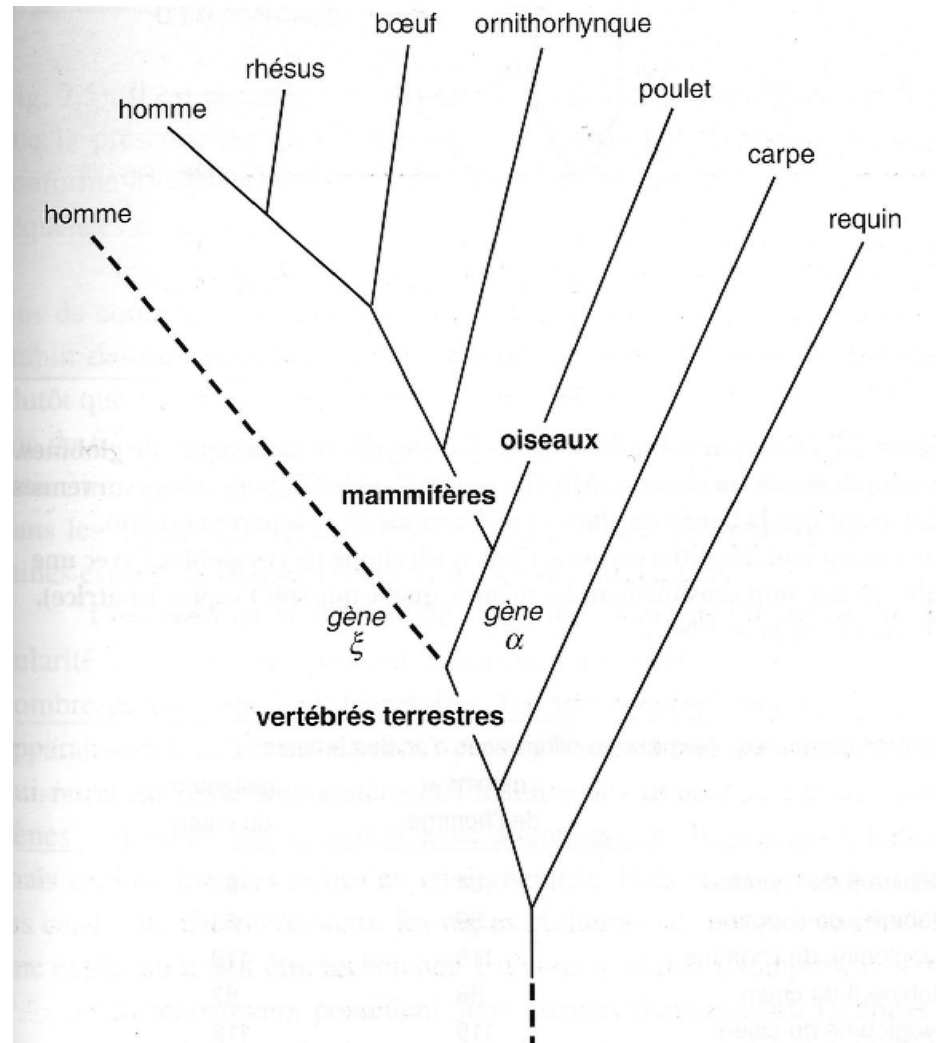


FIG. 2. Consensus maximum parsimony tree of IleRS, ValRS, and LeuRS genes using the program PROTPARS (26). Numbers are the frequency of occurrence of nodes in 300 bootstrap replicates.

## Duplications géniques & génomiques

### Identification des orthologues et des paralogues de globine



# Duplications géniques & génomiques

## Evolution des globines

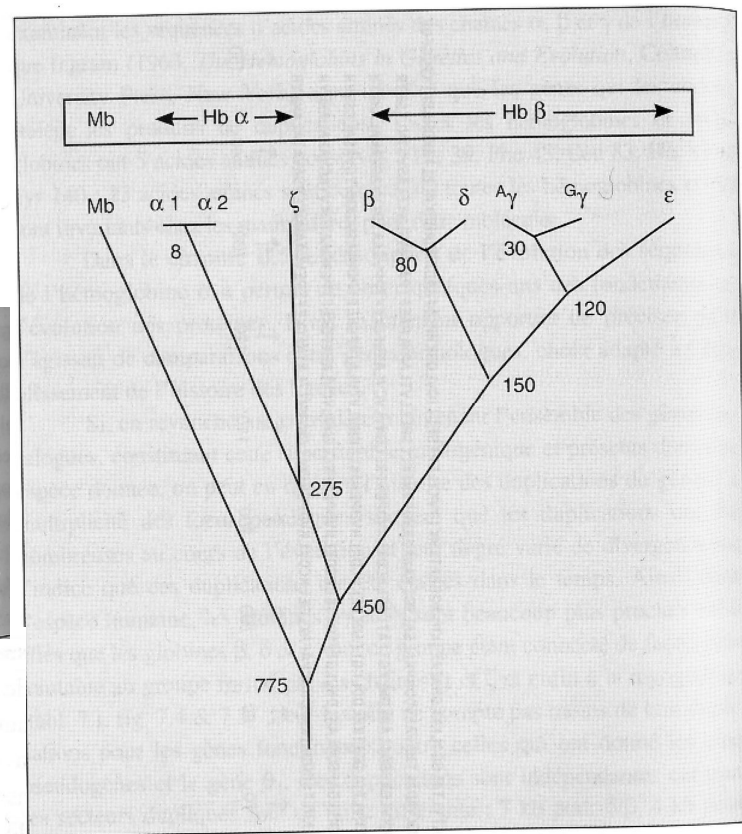
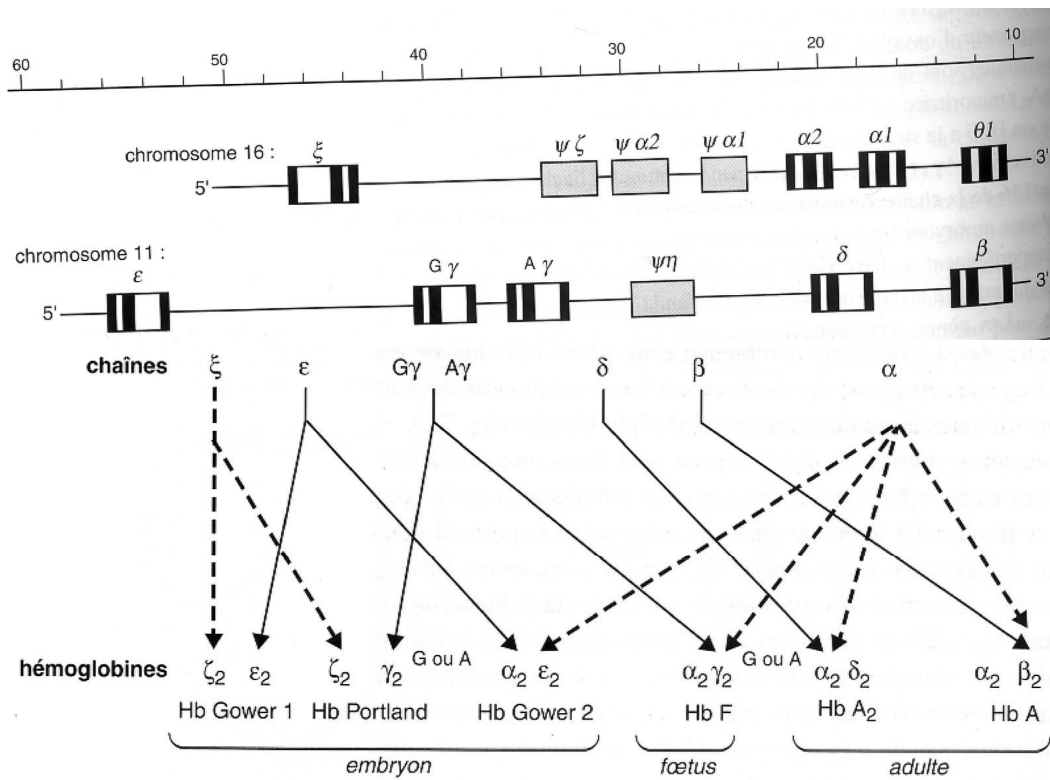
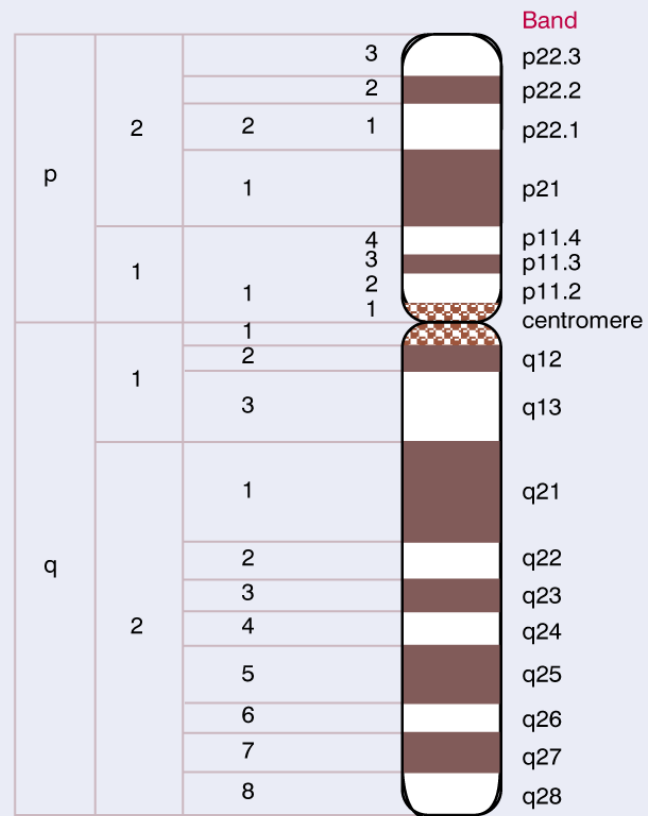


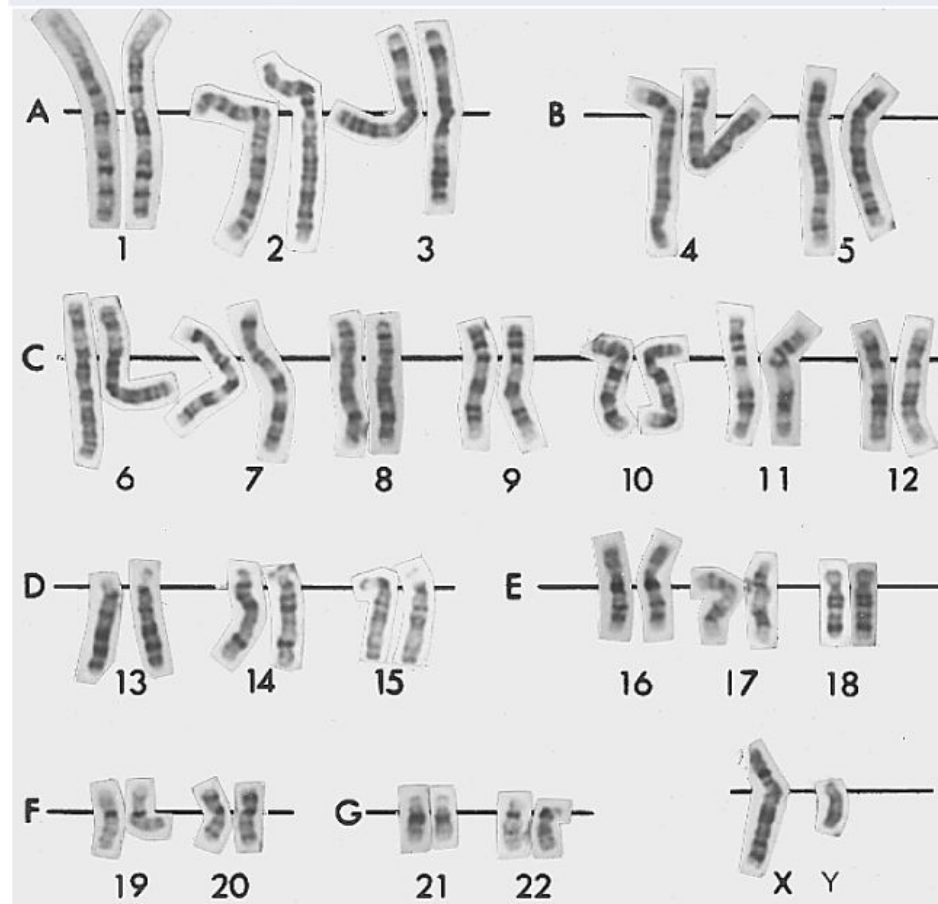
Figure 7.5 : Phylogénie des globines humaines.

L'âge approximatif des duplications est indiqué. Les duplications myoglobine/hémoglobine et Hb $\alpha$ /Hb $\beta$  sont très anciennes. La lignée  $\alpha$  comprend deux duplications : la dernière a donné  $\alpha_1$  et  $\alpha_2$ , il y a environ 8 MA ; les deux chaînes sont identiques. La diversification des gènes  $\beta$  comprend une première duplication, générant un proto- $\beta$  et un proto- $\epsilon$  vers 150 ou 200 MA. Les marsupiaux n'ont que ces deux gènes. À la base du rameau euthérien, vers 80-100 MA, la duplication de proto- $\beta$  a fourni  $\beta$  et  $\delta$  et la duplication de proto- $\epsilon$  a fourni  $\epsilon$ ,  $\gamma$  et  $\eta$ . Ce dernier est devenu un pseudogène chez les primates. Les chaînes  $A\gamma$  et  $G\gamma$  (duplication il y a 25 à 35 MA) diffèrent par 1 seul acide aminé (Ala ou Gly en position 136).

**Figure 18.12** The human X chromosome can be divided into distinct regions by its banding pattern. The short arm is *p* and the long arm is *q*; each arm is divided into larger regions that are further subdivided. This map shows a low resolution structure; at higher resolution, some bands are further subdivided into smaller bands and interbands, e.g. p21 is divided into p21.1, p21.2, and p21.3.

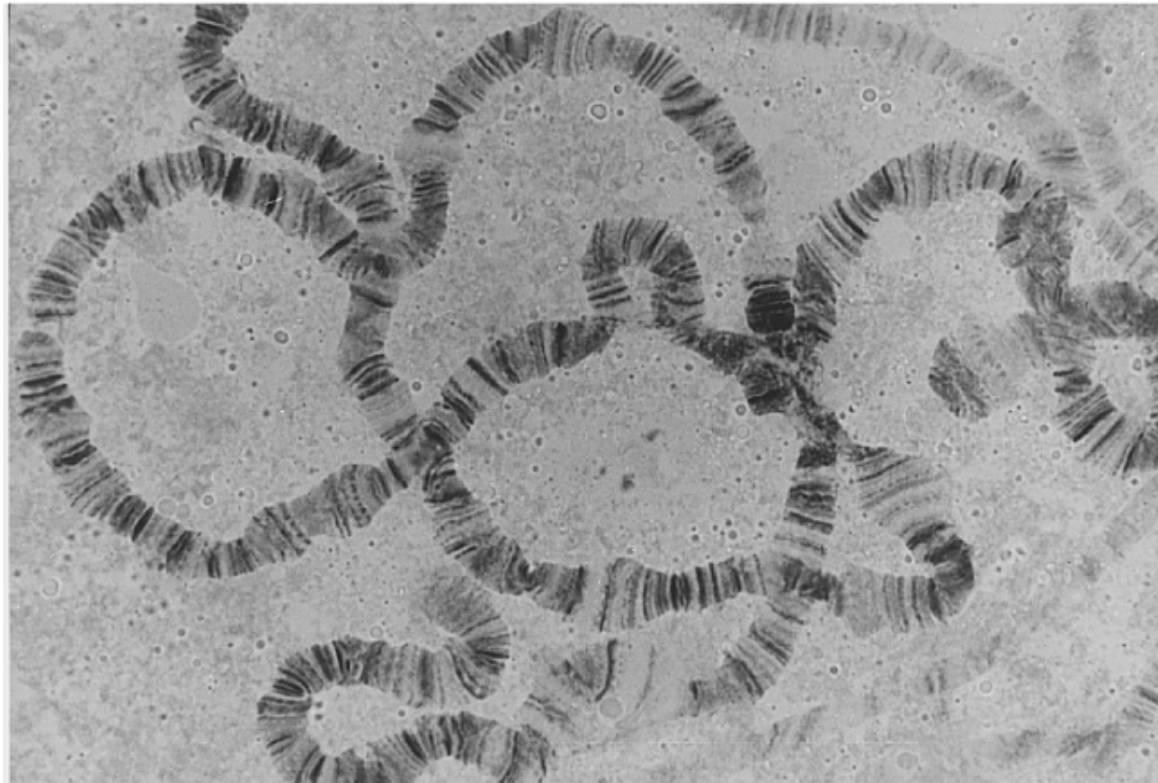


**Figure 18.11** G-banding generates a characteristic lateral series of bands in each member of the chromosome set. Photograph kindly provided by Lisa Shaffer.





**Figure 18.15** The polytene chromosomes of *D. melanogaster* form an alternating series of bands and interbands. Photograph kindly provided by Jose Bonner.

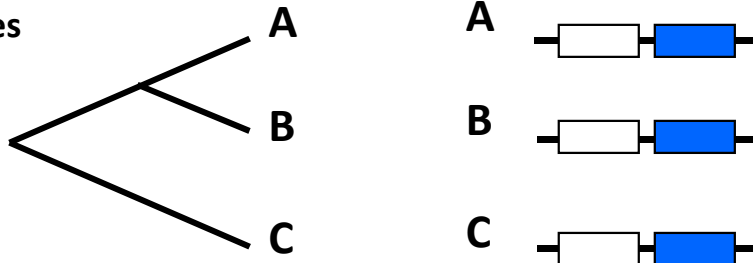


# Les bandes chromosomiques

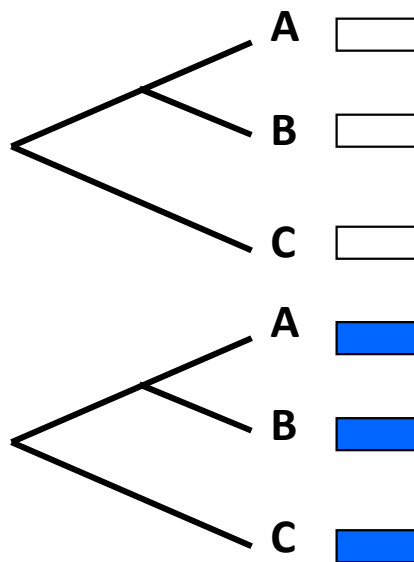
# Duplications géniques & génomiques

## Identification des orthologues et des paralogues

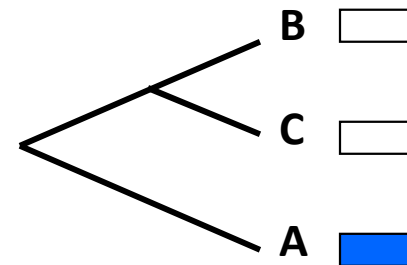
Les espèces



Phylogénies obtenues en comparant des orthologues



Phylogénie obtenues en comparant des orthologues et des paralogues  
(Si le paralogue chez A a fortement divergé)



Phylogénie d'orthologues = Phylogénie des espèces  
Phylogénie de paralogues = Phylogénie différente