

# Analyse de Séquences

M1 BIBS

V. 2019.1

# Contenu général de l'UE

- **D.Gautheret**
  - Les données de séquence: séquençage, structure des génomes, banques de données génomiques
- **O. Lespinet**
  - Alignement de séquence, phylogénie moléculaire
- **D.Gautheret – A. Lopes**
  - Homologie, annotation
  - TP: L'analyse de séquences en pratique: domaines, homologie, phylogénie

14h-17h		
<b>Jeudi 19/09/2019</b>	intro génomes / seq databases	DG
25-sept-19	Alignement séquences	OL
02-oct-19	Alignement séquences	OL
09-oct-19	Alignement séquences	OL
16-oct-19	Rech. Motifs	DG
23-oct-19	Arbres phylogénétiques	OL
	Vacances	
06-nov-19	exercices arbres	OL
13-nov-19	<b>Libre (Junior Conf)</b>	
<b>07-nov-19</b>	Annotation+TP1	DG
20-nov-19	TP2	DG
27-nov-19	TP3	DG
04-déc-19	TP4	AL
11-déc-19	prez par etudiants	AL
18-déc-19	TP5	DG+AL

# Préface

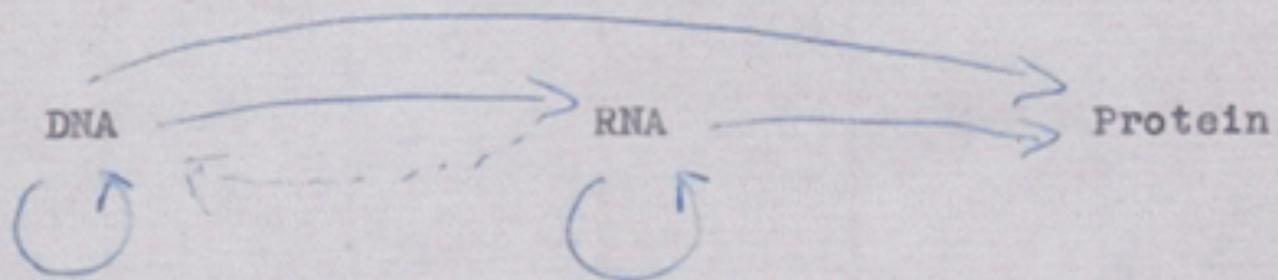
- Quelques notions d'information biologique

With slides from **Sacha Schutz**

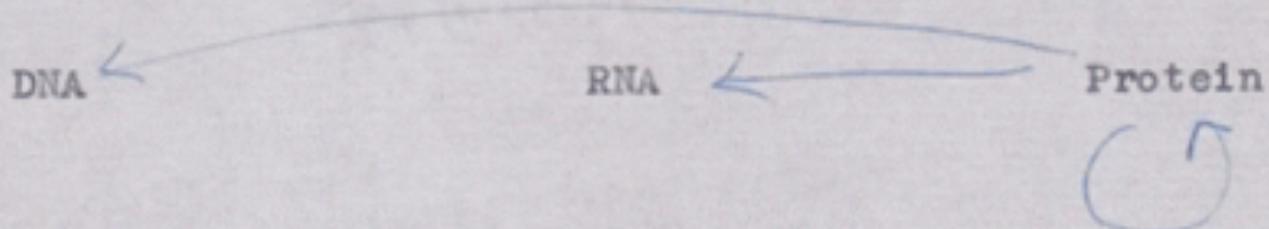
<https://github.com/dridk>

# Biologie=science de l'information?

The Central Dogma: "Once information has got into a protein it can't get out again". Information here means the sequence of the amino acid residues, or other sequences related to it. That is, we may be able to have

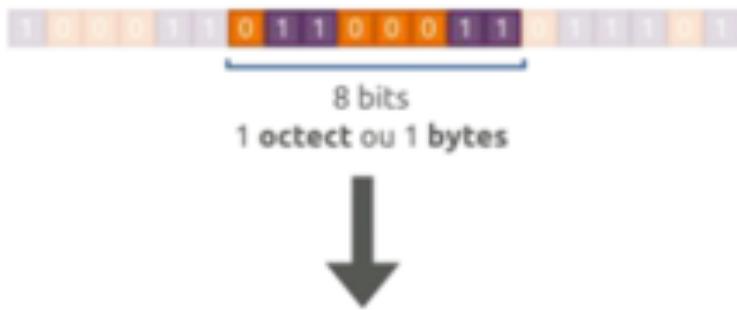
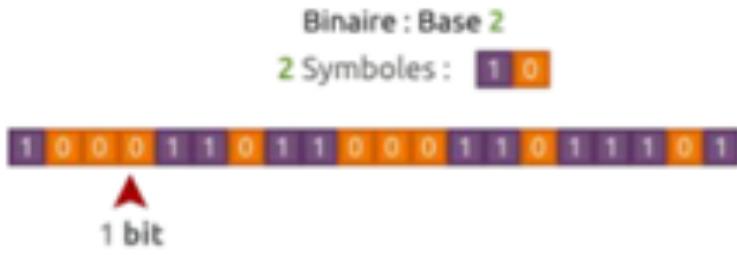


but never



where the arrows show the transfer of information.

# Biologie=science de l'information?



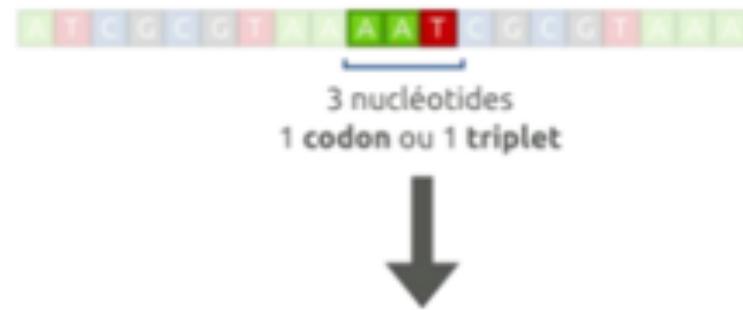
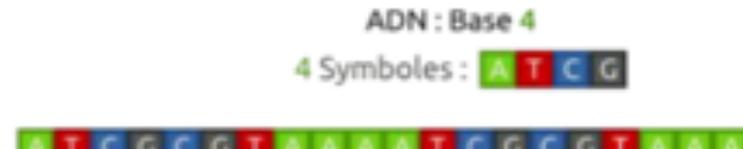
Combinaison = bases<sup>taille</sup> =  $2^8 = 256$

clef	valeur
1000001	A
1000010	B
1000000	⊕
0001101	<end line>

Encoder un texte

1000001	A
1000010	B
1000000	⊕
0001101	<end line>

• •  
Code ASCII



Combinaison = bases<sup>taille</sup> =  $4^3 = 64$

clef	valeur
ACG	Thr
AAG	Lys
GCG	Ala
TAG	<Stop>

Encoder une protéine

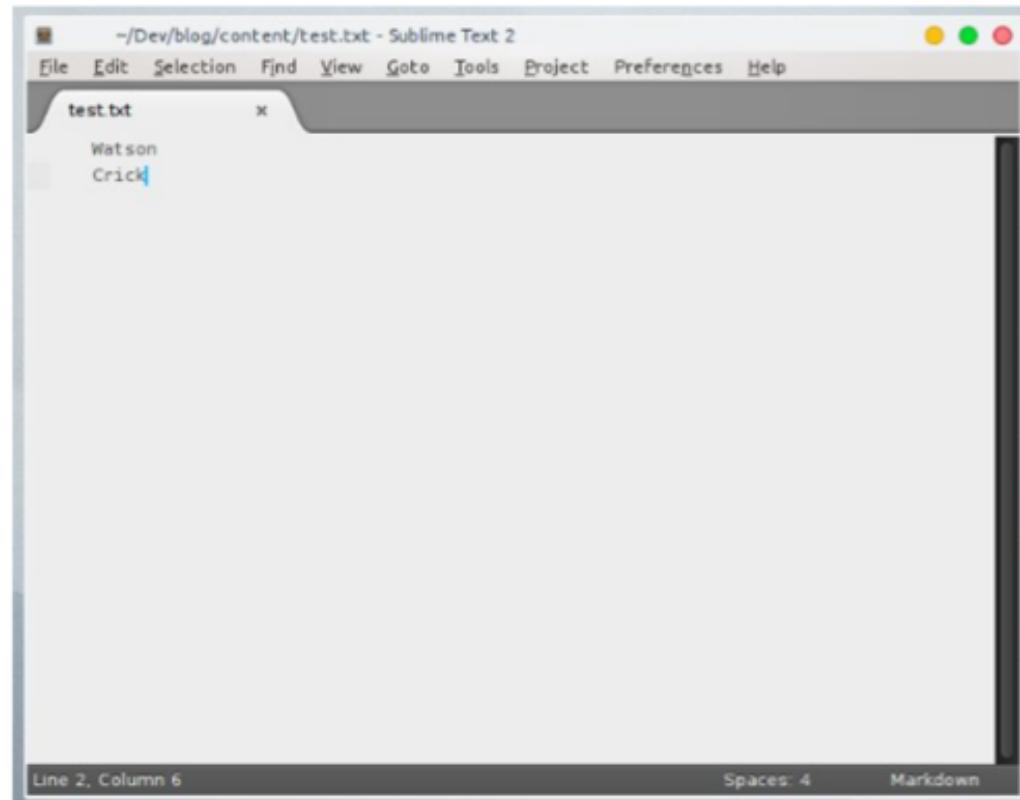
ACG	Thr
AAG	Lys
GCG	Ala
TAG	<Stop>

• •  
Code génétique

# Taille d'un fichier

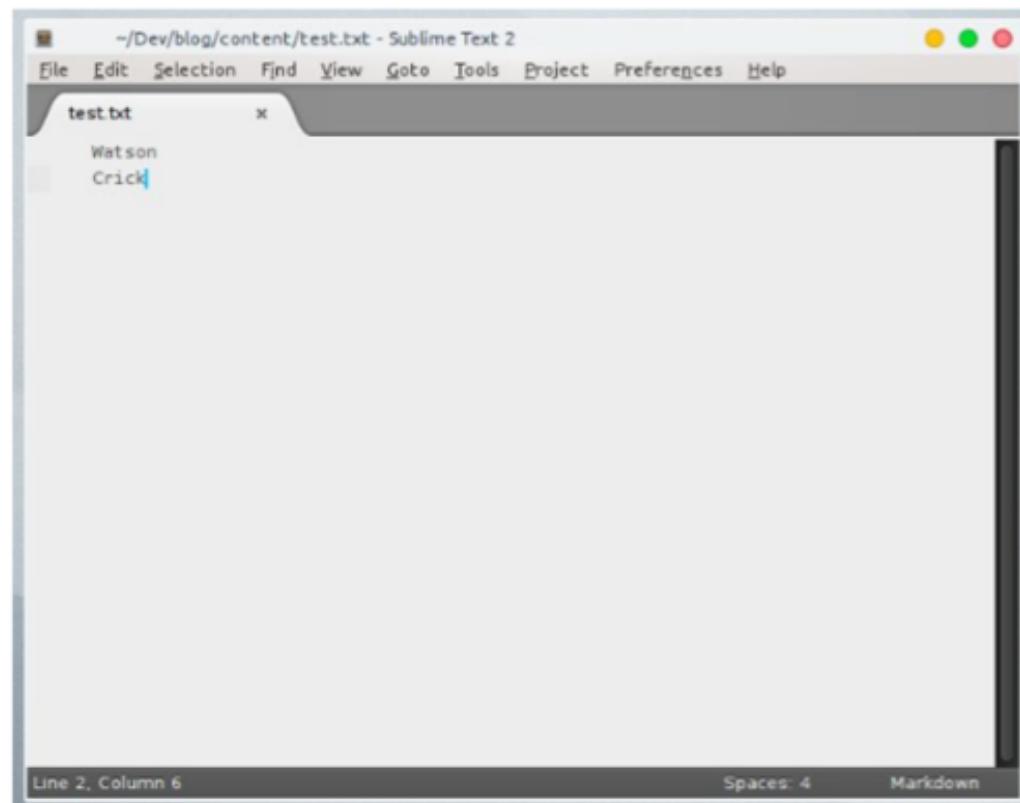
## Taille du fichier

- en octets?
- En bits?

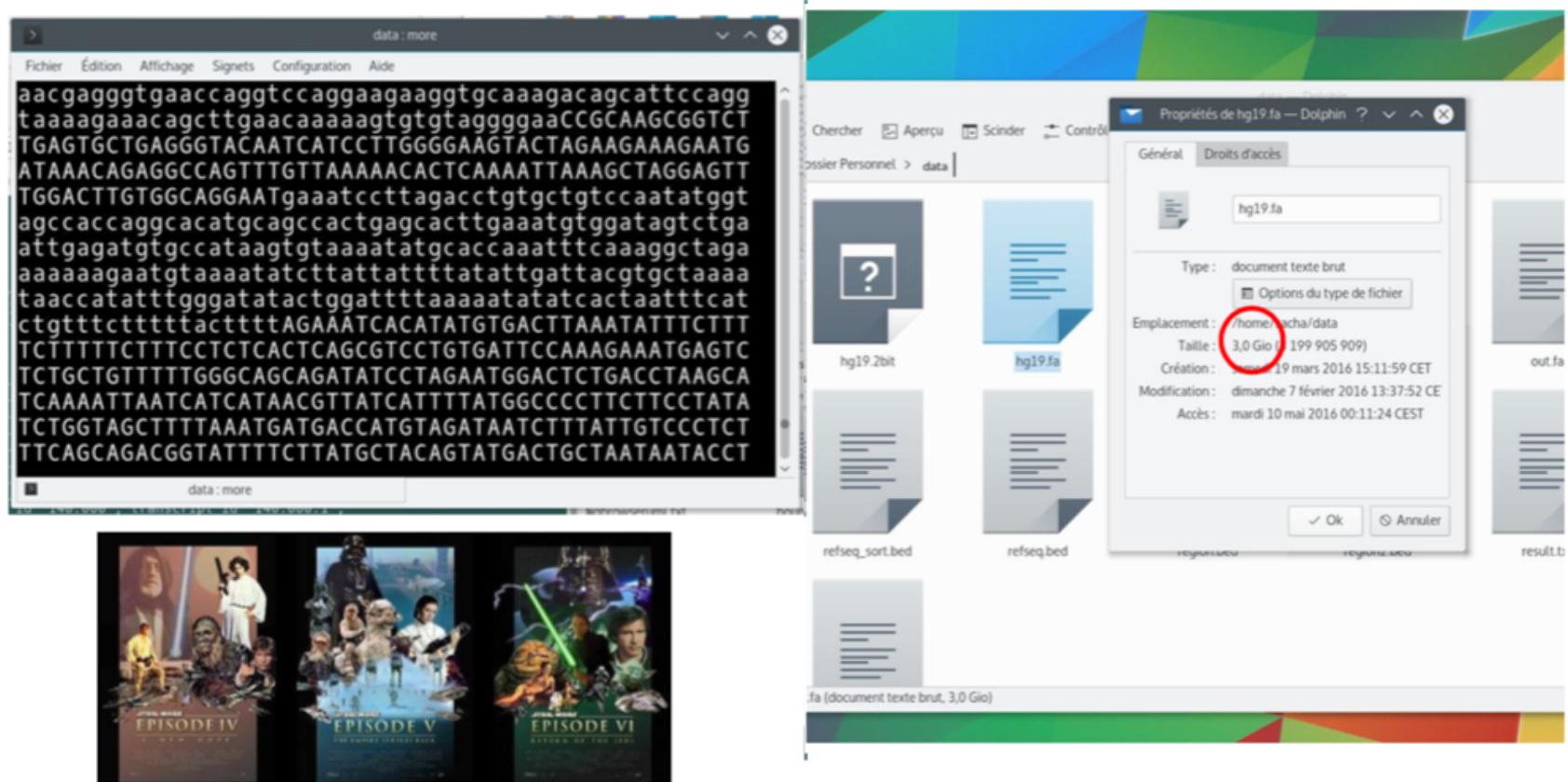


# Taille d'un fichier

12 octets  
(bytes)=11+1  
96 bits (12x8)



# Taille d'un fichier



# 2 types de fichiers

## Fichier texte

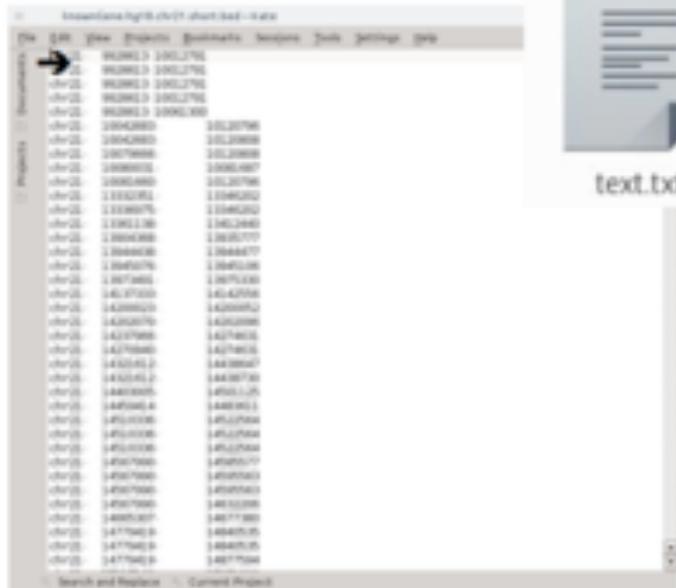
- Unité: 1 octet (byte)
- Lisible par un humain
- Dans un éditeur de texte
- Prend beaucoup d'espace
- Exemple:
  - HTML, Fasta, sam, csv, xml, vcf, fastq

## Fichier binaire

- Unité: 1 bit
- Non lisible par un humain
- Dans un éditeur hexadécimal
- Prend moins d'espace
- Exemple:
  - png, jpg, mp3, raw, bam, excel, word, fastq.gz

## Exemples

## Fichier texte



## Fichier binaire



# Espace texte vs. binaire

Exemple de fichier texte  
pour information  
Vrai/Faux (V/F):

VFVFFVFF

Total: 8 octets

Exemple de fichier  
binaire pour information  
Vrai/Faux (1/0):

10100100

Total: 1 octet

# Formats de données en bioinformatique

- La majorité des données de bioinfo sont de type texte:
  - FASTA, FASTQ, GB, SAM, VCF, BED, GFF, GTF, TSV, CSV, WML, JSON, PDB
- Pour des raisons de performance et d'espace, certains sont en format binaire
  - BAM, VCF.GZ, FASTQ.GZ

# Format et spécification

- Le format d'un fichier décrit comment les données sont représentées
- Cette description est fournie dans un document appelé **spécification**.

# Même donnée, différents formats

```
users : {  
first_name: "James", last_name:  
"Watson", birthday: "1928-04-06"  
}
```

Format **JSON**

<https://tools.ietf.org/html/rfc4627>

```
<users>  
<first_name>James</firstname>  
<last_name> Watson</last_name>  
<birthday>19280406</birthday>  
</users>
```

Format **XML**

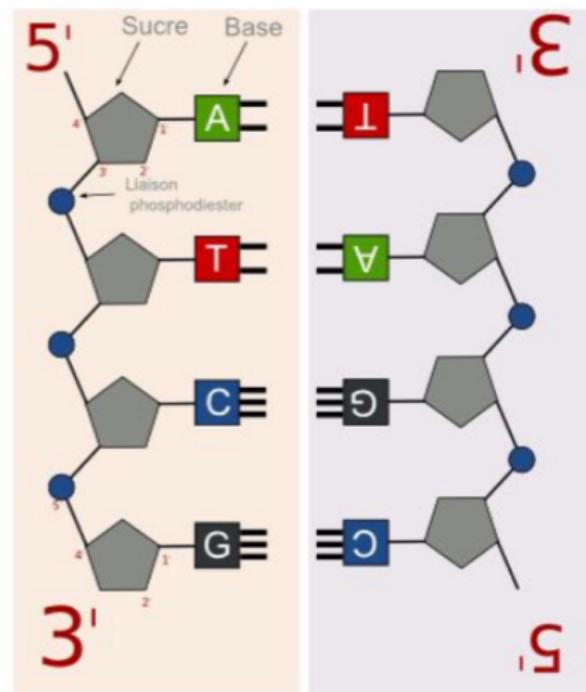
<https://www.w3.org/TR/REC-xml/>

# Séquences et régions

- En génomique on peut catégoriser les formats en:
  - Formats décrivant des séquences
  - Formats décrivant des régions

# Séquences d'ADN

- Toujours dans le sens 5'->3'
- Sur quel brin?



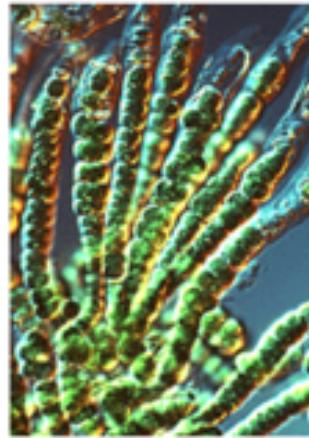
# Format fasta

\*.fa , \*.fasta

```
>identifiant1 commentaire libre
CAGCATCGATCGTCGGCGATGCATGCGGATGCTAGCTGATCACGATGC
CGCATGCTAGTCAGGCAGGGAGGGATATTATTAGCGGCGTATCGGATGA
CAGCATTACGGCGGGAGTGCTATTATTATGAGCGGGCGAT
>identifiant2 commentaire libre
CAGGCAGGAGGTTCTTATTATATCGGCAGGGCGGAGGCAGGCGATGCATC
CAGTGCAGTACGTAGTCAGCGATGCATTATGACTGACTCAGTTT
CCCGCTAGCTATGCTATGCTATTGATCGATTGAGCTGAGCTGATCTGGC
CAGCTATGCTTAGTA
```

# 1. Structure des génomes et des gènes

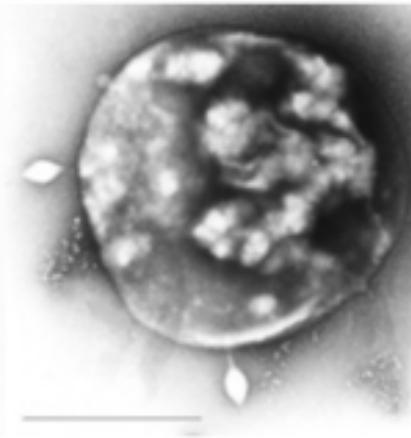
# The tree of Life



**B**



**E**



**A**

Picture: <http://tolweb.org/>

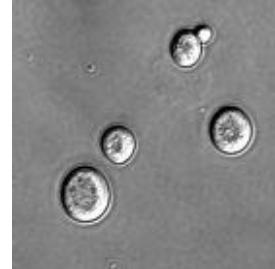
# Les organismes modèles

- **Facilité expérimentale**
  - Grosses cellules, corps transparent, temps de génération court, taux de reproduction élevé
- **Représentant d'une fonction**
  - Différenciation cellulaire, Système immunitaire, Photosynthèse, symétrie bilatérale...
- **Génome dense**
- **Données accumulées**
  - Mutants, génétique, biochimie...

# Les espèces-modèle



*E. Coli*  
4.5 Mb  
4200 gènes



*S. cerevisiae*  
13 Mb  
6000 gènes



*Dictyostelium*  
34 Mb  
12500 gènes



*Arabidopsis thaliana*  
150 Mb  
25000 gènes



*C. elegans*  
100 Mb  
15000 gènes



*Drosophila melanogaster*  
150 Mb  
15000 gènes



*Danio rerio*  
1,5 Gb  
45000 gènes



*Mus musculus*  
3 Gb  
20000 gènes

Et beaucoup d'autres ...

# *Homo sapiens*

- **Phylum**
  - Métazoaire, Mammifère, Primate
- **Inconvénients**
  - Expérimentation impossible sauf cellules somatiques en culture (Hela..)
  - Temps de génération lent



Raël

3 Gb  
20000 gènes codants

# Caractéristiques des génomes procaryotes (bactéries et archées)

- Chromosome circulaire unique
- Présence possible de petites séquences d'ADN circulaires indépendantes : **les plasmides.**

# Caractéristiques des génomes eucaryotes

- Noyau
- Taille >> procaryote
- Plusieurs chromosomes (homme 23, cheval 32, levure 16, drosophile 4...)
- La taille n'est pas proportionnelle à la complexité



Blé: 5.5 Gbases



3 Gbases

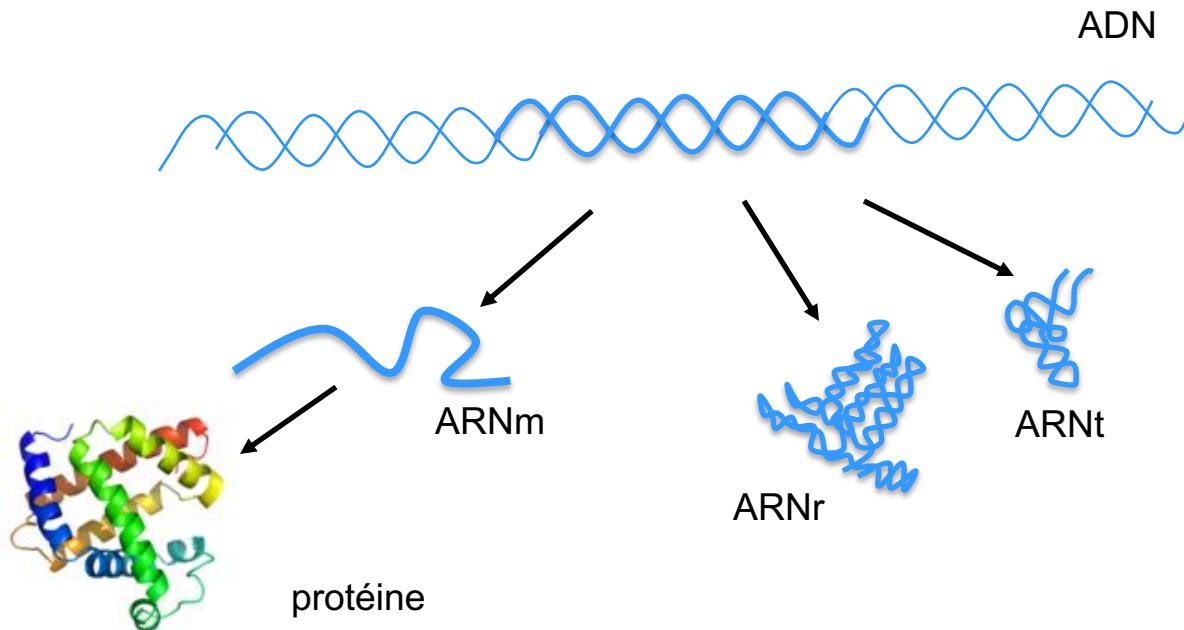
# Metagénomes

- Fragments d'ADN séquencés aléatoirement à partir d'un environnement donné
- Soit fragments ciblés (p. ex. ARNr 16S), soit ADN génomique
- Environnements étudiés:
  - Océan, intestin humain, drainage minier acide,
- Applications:
  - Populations microbiennes (structure, variations)
  - Découverte et séquençage de microorganismes
  - Découverte de gènes/fonctions

# Séquence des gènes

# Le gène

- Un gène est une séquence d'ADN qui spécifie la synthèse d'une protéine ou d'un ARN fonctionnel
- Un gène peut donc coder pour un ARN messager ou pour un ARN non-messager (ARNr, ARNt, ...)

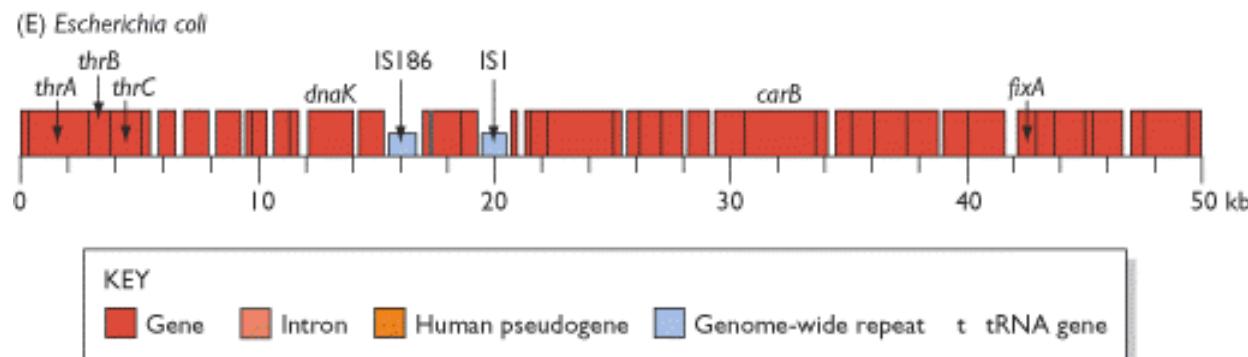


# Les gènes procaryotes

- Fraction codante des génomes élevée.
  - > 90% codant
  - Peu de séquences intergéniques
  - Génome « compact »
- Chez les procaryotes **la séquence des gènes est continue. Pas d'intron**
- Gènes organisés en opérons. 600 opérons dans le génome de *Escherichia coli*.

# Densité des gènes procaryotes

- Longueur gène 950 nt. en moyenne (coli)
- Haute densité en gènes: 1gène / kb
- 95% du génome est transcrit chez E. coli.

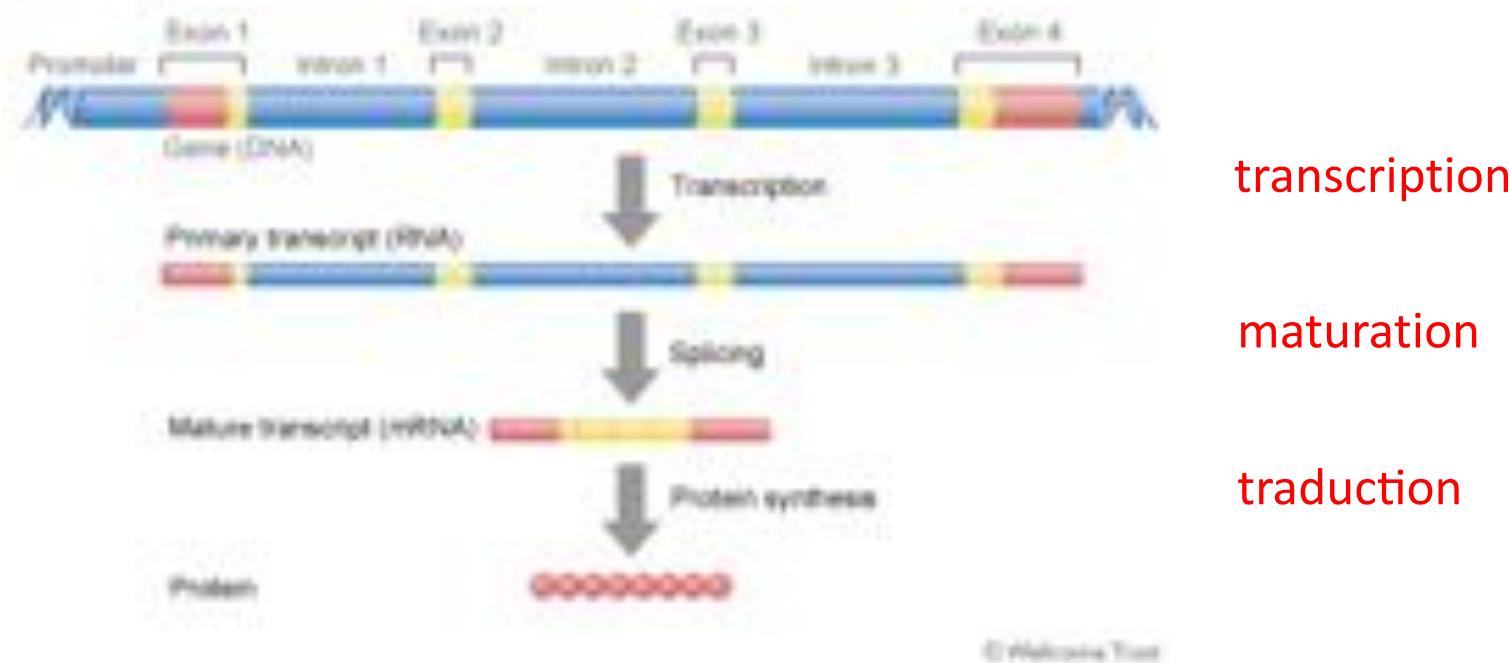


From « Genomes 2 », T.A. Brown

## Les gènes eucaryotes

- Gènes « disloqués » (exons, introns)
- Grandes régions intergéniques de fonction inconnue

# Le gène eucaryote



- Gène humain moyen: 27kb, 9 introns, codant:1,3kb , exon moyen: 145 bp, intron moyen:3365 bp.
- Gènes "monstres": dystrophine: 2,4 Mb; Facteur de coagulation VIII: 186 kb, 26 exons; Tinine: codant: 80kb, 178 exons

# Densité des gènes eucaryotes

Densité moyenne:

– *S. Cerevisiae*:

1 gène/2kb.



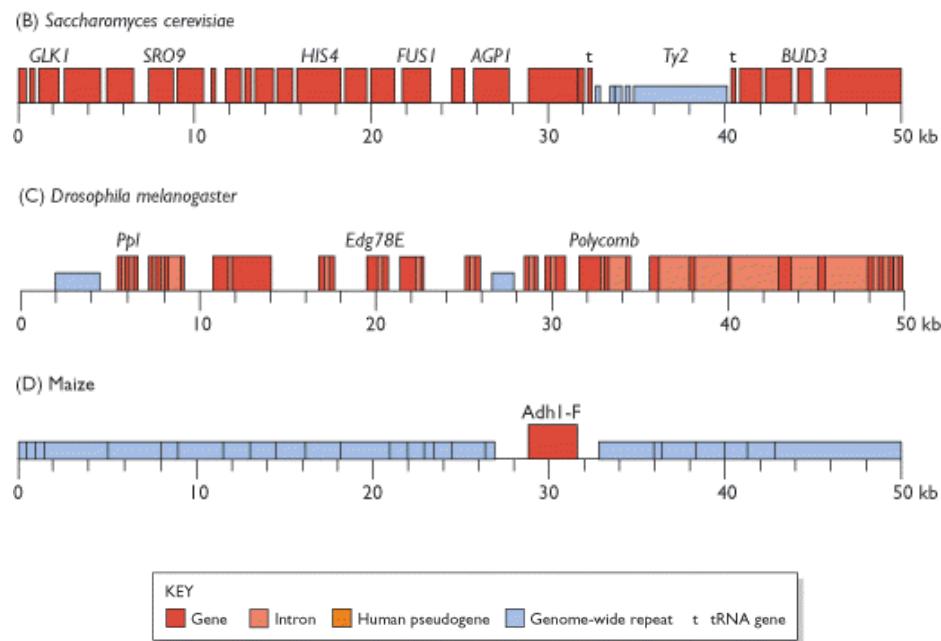
– *Drosophila*:

1 gène/10kb



– Maïs: 1 gène tous les 70kb

– Humain: 1 gène tous les 100kb



From « Genomes 2 », T.A. Brown

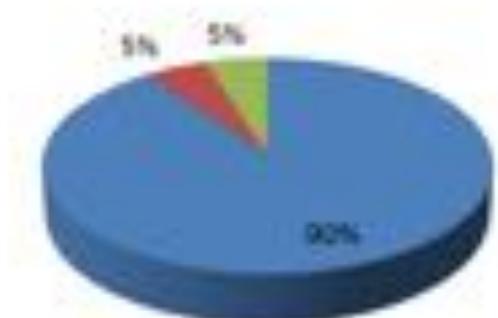
# Junk DNA: les séquences répétées dans le génome humain

- **4 classes de séquences répétées**

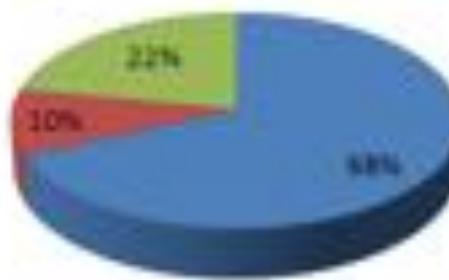
- Répétition de type transposon (ou interspersed repeats)
- copie inactives de gènes (processed pseudogenes)
- Répétition simples de k-mères courts, p. ex. (A) $n$ , (CA) $n$  ou (CGG) $n$
- Segments dupliqués: blocs de 10–300 kb copiés d'une région à l'autre ou en tandems



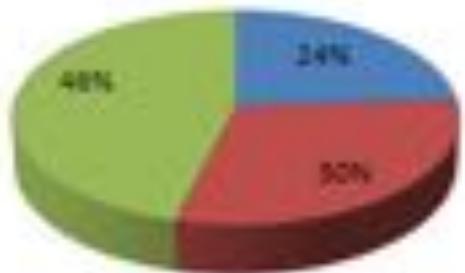
# Fraction codante et non-codante des génomes



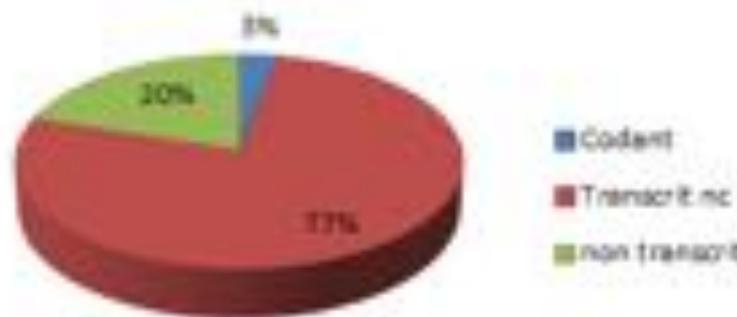
*E. coli*



*S. cerevisiae*



*C. Elegans*



*H. sapiens*

# Le format GTF (un format de régions)

- GTF: general feature format

1. **seqname** - The name of the sequence (chromosome/scaffold)
2. **source** - The program that generated this feature
3. **feature** - Type of feature ("CDS", "start\_codon", "stop\_codon", "exon")
4. **start** - Starting position of the feature in the sequence (starts at 1)
5. **end** - Ending position of the feature (inclusive).
6. **score** - Score between 0 and 1000 (or ":" if no value)
7. **strand** - '+', '-' or ':'
8. **frame** - If coding exon, *frame* should be 0-2: reading frame of the first base.
9. **group** - All lines with the same group are linked together into a single item.

# Le format GTF

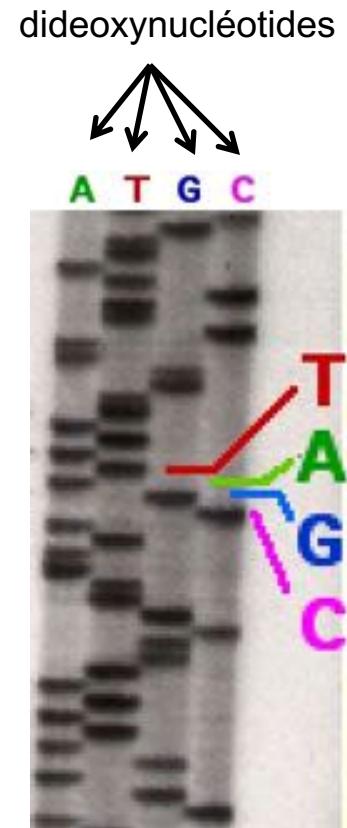
chr9	hg38_refGene	stop_codon	133255666	133255668	0.000000	- .	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	CDS	133255669	133256356	0.000000	- 1	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	exon	133255176	133256356	0.000000	- .	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	CDS	133257409	133257542	0.000000	- 1	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	exon	133257409	133257542	0.000000	- .	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	CDS	133258097	133258132	0.000000	- 1	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	exon	133258097	133258132	0.000000	- .	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	CDS	133259819	133259866	0.000000	- 1	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	exon	133259819	133259866	0.000000	- .	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	CDS	133261318	133261374	0.000000	- 1	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	exon	133261318	133261374	0.000000	- .	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	CDS	133262099	133262168	0.000000	- 2	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	exon	133262099	133262168	0.000000	- .	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	CDS	133275162	133275189	0.000000	- 0	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	start_codon	133275187	133275189	0.000000	- .	gene_id	NM_020469;	transcript_id	NM_020469;
chr9	hg38_refGene	exon	133275162	133275214	0.000000	- .	gene_id	NM_020469;	transcript_id	NM_020469;

- Avec
  - un GTF
  - un FASTA
- on peut reconstruire le transcriptome entier d'un organisme

## 2. Le séquençage

# Le séquençage de Sanger (1977)

- Séquençage par terminaison de chaîne
  - Utilisation de dideoxynucléotides pour interrompre la synthèse à un certain type de base.
  - 4 réactions + marquage radioactif
- Amélioré en 1987 par l'introduction de marqueurs fluorescents (1 seule réaction) et l'automatisation.



Wikipedia

# Stratégie de séquençage Shotgun

Clone à séquencer

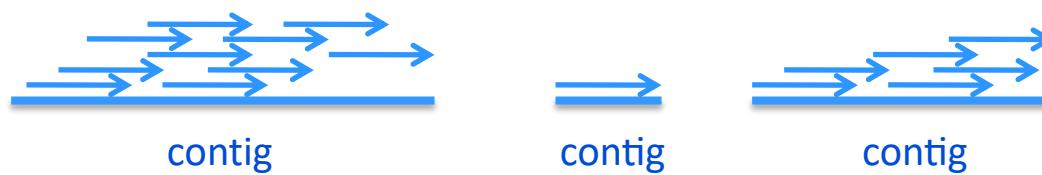
Fractionnement



Séquençage aléatoire



Assemblage



Scaffold (si information cartographie)



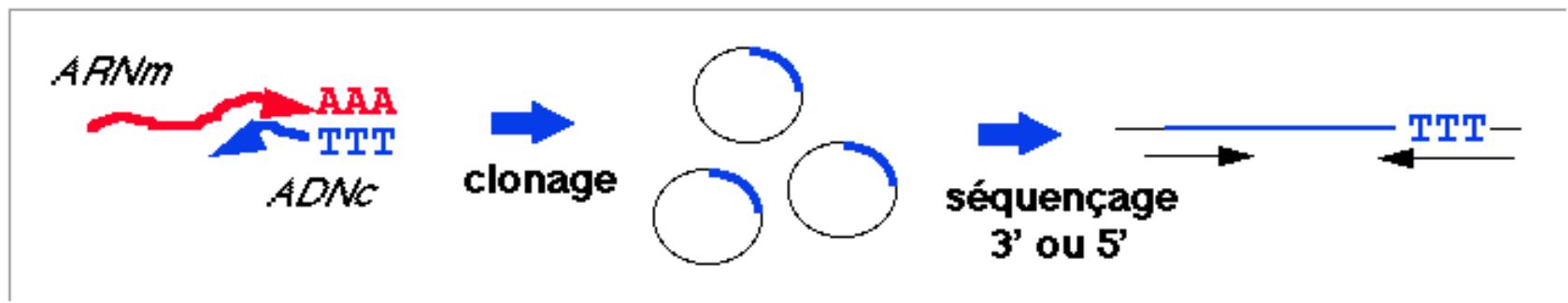
Problématique  
en présence de  
longues régions  
répétées)



Eric Lander, coordinateur: projet de séquençage du génome humain

# Les cDNA

- Idée initiale:
  - pourquoi vouloir tout séquencer (95% de junk DNA) si ce sont les gènes qui nous intéressent?
- Usage actuel:
  - Analyse du transcriptome

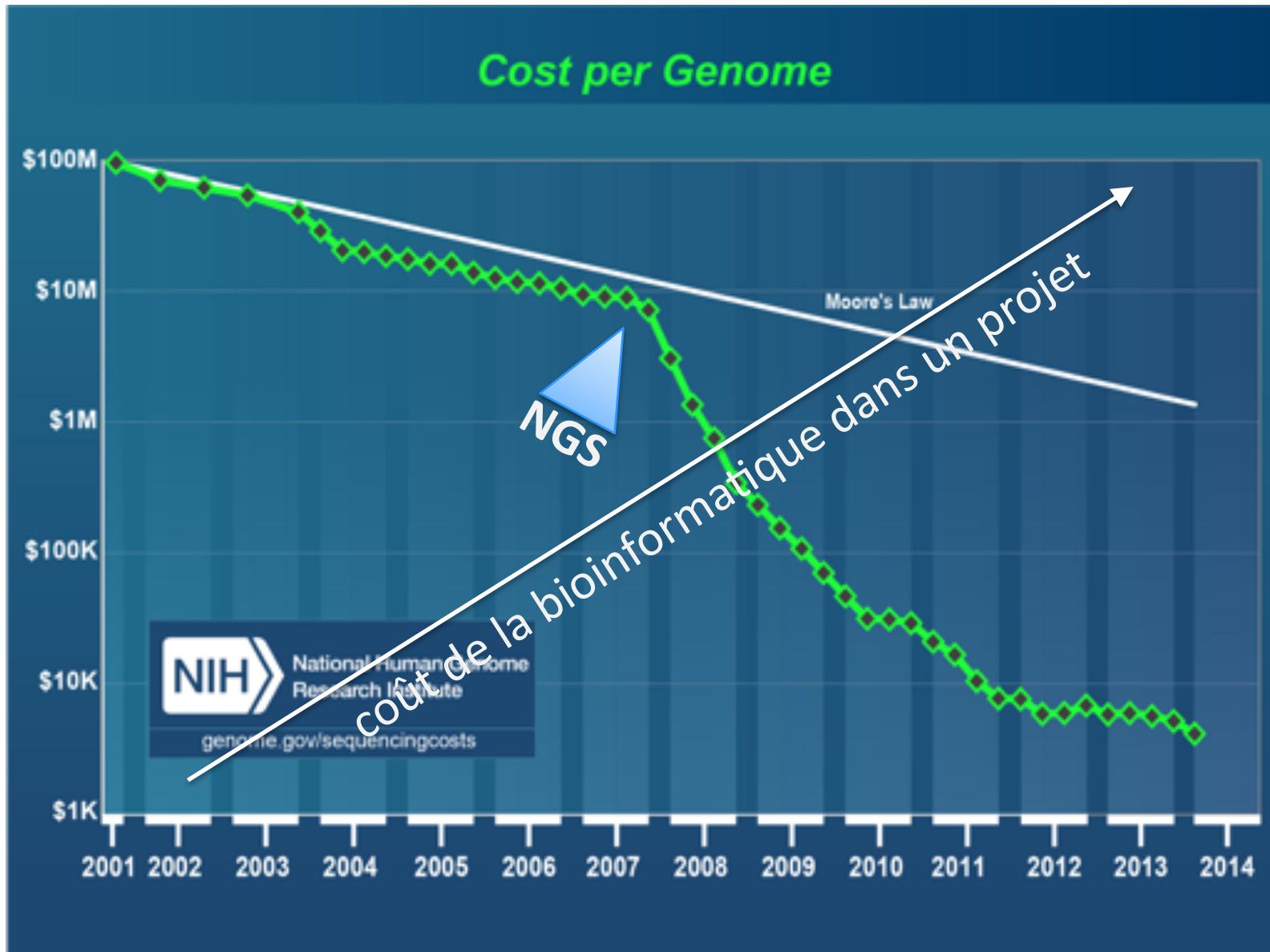


- EST (Expressed sequence Tag) = Séquences partielles d'ADNc clonés et prélevés aléatoirement.
- Full length cDNA: Séquences complète d'ADNc



Craig Venter a contribué à l'usage massif des EST dès 1991

# Le bouleversement des NGS

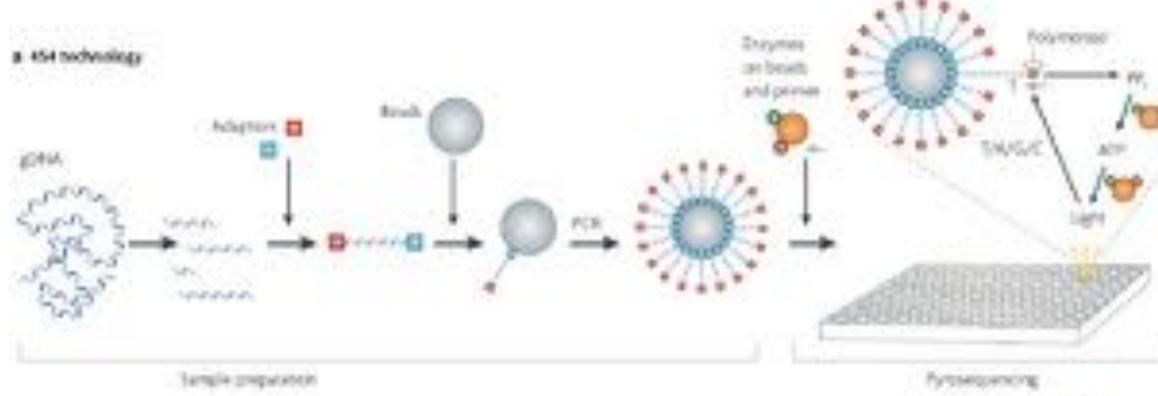
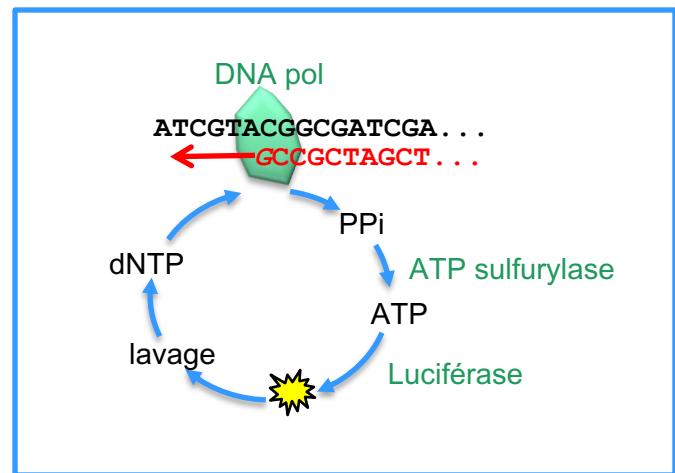


# Début des NGS: pyroséquençage (2006)

454 LifeSciences/Roche  
Biotage/Qiagen

1. ADN immobilisé (billes)
2. Synthèse brin complémentaire par ADN polymérase
3. Introduction des dNTP un par un
4. Si bon dNTP: libération PPi, synthèse ATP, et émission de lumière.

300-500nt à la fois x 400.000



Nature Reviews Microbiology

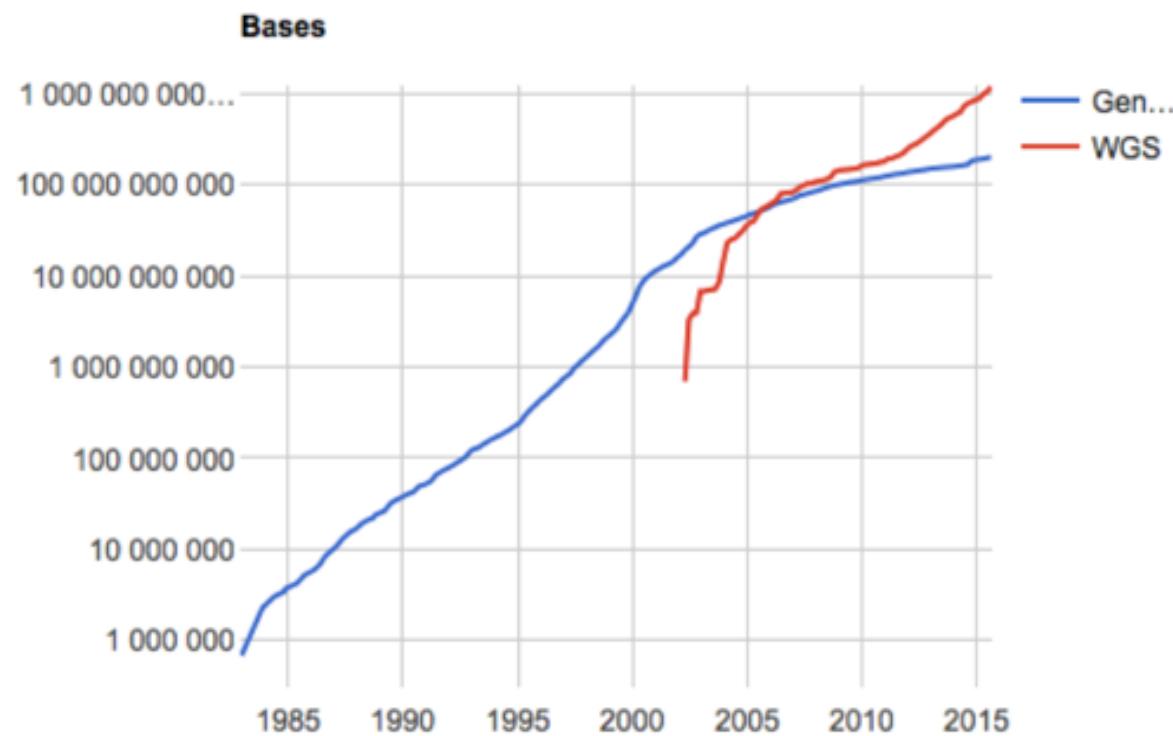
### 3. Les banques de données génomiques

# Genbank: La banque d'ADN du NIH

## Uniquement des séquences assemblées

- **Etat 2018**

- $> 10^{11}$  bases
- $> 10^8$  séquences
- Genbank double environ tous les 14 mois depuis ses débuts en 1982.
- Nouvelle version tous les 2 mois



# Enregistrement Genbank

- Chaque enregistrement se voit attribuer un numéro d'acquisition, stable et unique, et chaque séquence un numéro GI.
- Quand un changement est effectué dans un enregistrement Genbank, le num. d'acquisition reste, le GI change.

# Séquence au Format Fasta

```
>U00096 Escherichia coli K-12 MG1655 complete genome.  
agctttcattctgactgcaacgggcaatatgtctctgtgtggattaaaaaaaagagtgtc  
tgatagcagttctgaactggttacctgccgtgagtaaattaaattttattgacttagg  
tcactaaatactttaaccaaatataggcatagcgcacagacagataaaaattacagagtac  
acaacatccatgaaacgcattagcaccaccattaccaccaccattaccacaggt  
aacgggtgcgggctgacgcgtacaggaaacacagaaaaaagccgcacctgacagtgcggg  
cttttttttcgaccaaaggtaacgaggttaacaaccatgcgagtgtaagttcggcggt  
acatcagtggcaaattgcagaacgtttctgcgtttgcgttatattctggaaagcaatgcc  
aggcaggggcaggtggccaccgtcctctgcggccggccaaatcaccaaccacctggtg  
gcgatgattgaaaaaccattagcggccaggatgcttacccaatattcagcgatgccgaa  
cgtattttgccgaactttgacggactcgccgcgcgcggccagccgggttcccgtggcg  
caattgaaaacttcgtcgatcaggaatttgcggccaaataaaacatgtcctgcattggcatt  
agttgttggggcagtgcggatagcatcaacgcgtgcgtgatattgcgtggcgagaaa  
atgtcgatcgccattatggccggcgtattagaagcgcgccgtcacaacgttactgttac  
gtccggcgtggaaaaactgctggcagtggggcattacctcgaaatctaccgtcgatattgct  
gagttccacccggcgtattgcggcaagccgcattccggctgatcacatgggtctgatggca  
ggttcaccggccggtaatgaaaaaggcgaactgggtgctggacgcaacgggtccgac  
tactctgctgcggcgtggctgcctgtttacgcggcattgtgcgagattggacggac
```

# Enregistrement Genbank avec annotation

LOCUS L10986 47233 bp DNA linear INV 21-SEP-2004  
 DEFINITION Caenorhabditis elegans cosmid F10E9, complete sequence.  
 ACCESSION L10986  
 VERSION L10986.2 GI:38638818  
 KEYWORDS HTG.  
 SOURCE Caenorhabditis elegans  
 ORGANISM Caenorhabditis elegans  
     Eukaryota; Metazoa; Nematoda; Chromadorea; Rhabditida;  
     Rhabditoidea; Rhabditidae; Peloderinae; Caenorhabditis.  
 REFERENCE 1 (bases 1 to 47233)  
 AUTHORS .  
 CONSRTM WormBase Consortium  
 TITLE Genome sequence of the nematode *C. elegans*: a platform for investigating biology. The *C. elegans* Sequencing Consortium  
 JOURNAL Science 282 (5396), 2012-2018 (1998)  
 MEDLINE 99069613  
 PUBMED 9851916  
 FEATURES Location/Qualifiers  
     source 1..47233  
           /organism="Caenorhabditis elegans"  
           /mol\_type="genomic DNA"  
           /strain="Bristol N2"  
           /db\_xref="taxon:6239"  
           /chromosome="III"  
           /clone="F10E9"  
     gene 265..26728  
           /gene="mig-10"  
           /locus\_tag="F10E9.6"  
     CDS join(265..338,3266..3515,15194..15317,21507..21  
           21727..21887,23171..23335,24302..24472,24524..24608,  
           25012..25827,26284..26430,26478..26728)  
           /gene="mig-10"  
 /translation="MDSCEECDLEVDSDEEDQLFGEKCISLSSLLPLSSSTLLSNA"

BASE COUNT 2598 a 2024 c 1888 g 2449 t  
ORIGIN  
1 ttctaaaagt cgaaaaacga gcaattttt atgctagatt ttttgatttg acgaattttt  
61 tcagttttt ttctttaaaa aaggttttt accccttaaa gtttccttt ccctcccaat  
121 ttttcccttc ttctttatac gacttctcaa gtttcaactc taaaacaag ctacatgtac  
181 atttccggta aactttgtgt ctcagaagat ccattttctt tttgttacat ttattcaaga  
241 ttgaattcca aaatttcagc caatatggac agttgcgaag aggaatgcga tctggaagtt  
301 gacagtacg aagaagatca actttttgtt gaaaagtggt gagtttctt tggtttaacc  
361 aaagaatcatcgt cagttgtccg taaacacttg actcccaaat ggtttctcgtaattaccta  
421 tgcaacttcttcaacttgcgtt gccgtttgtt cttagccaat ttggaaacgtt tagatgttaa  
481 atggaaaatgg tggtaaagg tttatttttt agaaaaaaaaagg tttggaaaaaa aatcgagtca  
541 ctgaatagg tgaagaacgg aaaataaaa ctttccaaaa atcataaaaac atttagtgg  
601 tcgaaaatgg tagtgtttt tttgttggta tggttttgaca aaagctaaac catctttatt  
661 gtatgtttgtt aaaaatgttca caaagatgcg ttttttttttcaaaatggca ggctatctt  
721 acatccatcat ttggataataat tttatggtaa ttatcgctaa caaattttc tattttccca  
781 attttatcggtt ttatggtaa aacgg tttgttggta tttttgttctt atctttatgtt qtcatcgat

# Mini TD

- Sur le site NCBI/genome, retrouvez le génome annoté de *Candidatus Carsonella ruddii* HT.
- De quel type d'espèce s'agit-il?
- Taille du génome?
- Que contiennent les champs...
  - CDS
  - Complement
  - Translation
- Sauvez le fichier au format Genbank « full » (annotations et séquences) et comptez les CDS avec la commande grep, option -c
- A l'aide de wget, récupérez la séquence au format fasta et comparez la taille des fichiers fasta et gb. Récupérez le fasta avec wget (curl sur Mac)  
XXX=accession  
<https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=nucleotide&id=XXX&rettype=fasta&retmode=txt>

# Autres banques nucléotidiques

- EMBL: Equivalent européen de Genbank. Format différent, contenu presque identique.
- DDBJ: équivalent au Japon
- Banques spécialisées Certaines collections de séquences, bien que généralement présentes dans Genbank, sont beaucoup plus utiles lorsqu'elles sont rassemblées dans des banques spécialisées, par ex:
  - Récepteurs des lymphocytes T (Réarrangements de l'ADN)
  - Génomes HIV, etc.
- Banques pour Blast
  - NR nucléique (« Non-redundant »). All GenBank+EMBL+DDBJ+PDB sequences (but no EST, STS, GSS, or phase 0, 1 or 2 HTGS sequences). (n'est plus "non-redondant")
  - DbEST: dbest Database of GenBank+EMBL+DDBJ sequences from EST Divisions

# SRA: short read archive

- Partie de *European Nucleotide Archive* (sequences brutes issues de séquenceurs).
- Information sur chaque entrée:
  - Study
  - Sample
  - Experiment
  - Run
  - Organism
  - Instrument Platform
  - Library Name
  - Read Count
  - Base Count
  - File Name / File Size
- Séquences au format fastq

# Format fastq

Descriptif du read (position sur la piste de séquençage, taille,...)

```
@SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
GGGTGATGGCCGCTGCCGATGGCGTCAAATCCCACC
+SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
IIIIIIIIIIIIIIIIIIIIIIIIIIIIII9IG9IC
```



Qualité (probabilité que la base soit correcte) encodé par code ASCII

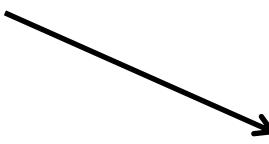
Attention: plusieurs versions (Illumina, Sanger..)

$$Q_{\text{sanger}} = -10 \log_{10} p$$

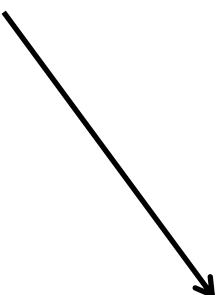
# Code ASCII et score de qualité

Décimal	Octal	Hex	Binaire	Caractère
000	000	00	00000000	NUL (Null char.)
001	001	01	00000001	SOH (Start of Header)
002	002	02	00000010	STX (Start of Text)
003	003	03	00000011	ETX (End of Text)
004	004	04	00000100	ETB (End of Transmission)
005	005	05	00000101	ENQ (Enquiry)
006	006	06	00000110	ACK (Acknowledgment)
007	007	07	00000111	BEL (Bell)
008	010	08	00001000	BS (Backspace)
009	011	09	00001001	HT (Horizontal Tab)
010	012	0A	00001010	LF (Line Feed)
011	013	0B	00001011	VT (Vertical Tab)
012	014	0C	00001100	FF (Form Feed)
013	015	0D	00001101	CR (Carriage Return)
014	016	0E	00001110	SO (Shift Out)
015	017	0F	00001111	SI (Shift In)
016	020	10	00010000	DLE (Data Link Escape)
017	021	11	00010001	DC1 (XON)(Device Control 1)
018	022	12	00010010	DC2 (Device Control 2)
019	023	13	00010011	DC3 (XOFF)(Device Control 3)
020	024	14	00010100	DC4 (Device Control 4)
021	025	15	00010101	NAK (Negative Acknowledgement)
022	026	16	00010110	SYN (Synchronous Idle)
023	027	17	00010111	ETB (End of Trans. Block)
024	030	18	00011000	CAN (Cancel)
025	031	19	00011001	EM (End of Medium)
026	032	1A	00011010	SUB (Substitute)
027	033	1B	00011011	ESC (Escape)
028	034	1C	00011100	FS (File Separator)
029	035	1D	00011101	GS (Group Separator)
030	036	1E	00011110	RS (Request to Send)(Record Separator)
031	037	1F	00011111	US (Unit Separator)
032	040	20	00100000	SP (Space)
033	041	21	00100001	! (exclamation mark)
034	042	22	00100010	" (double quote)
035	043	23	00100011	# (number sign)
036	044	24	00100100	\$ (dollar sign)
037	045	25	00100101	% (percent)
038	046	26	00100110	& (ampersand)
039	047	27	00100111	' (single quote)
040	050	28	00101000	( (left opening parenthesis)
041	051	29	00101001	) (right closing parenthesis)
042	052	2A	00101010	* (asterisk)
043	053	2B	00101011	+ (plus)
044	054	2C	00101100	, (comma)
045	055	2D	00101101	- (minus or dash)
046	056	2E	00101110	. (dot)
047	057	2F	00101111	/ (forward slash)
048	060	30	00110000	0
049	061	31	00110001	1
050	062	32	00110010	2
051	063	33	00110011	3

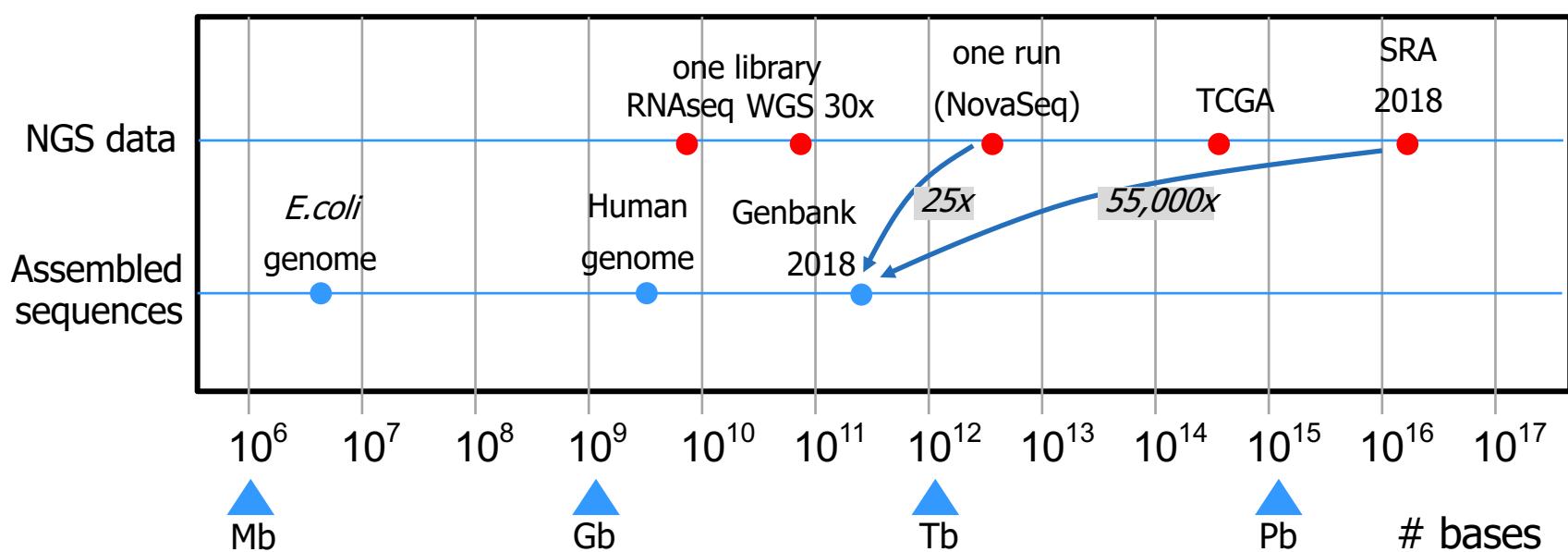
Score le plus bas  
(33=!): probabilité  
d'erreur élevée



Score le plus haut  
(126=~/): probabilité  
d'erreur faible



# The incredible scale of NGS databases



# Banques protéiques

- Swissprot (UniProtKB/Swissprot).
  - La mieux annotée des banques protéiques. 2011: 530.000 entrées.
  - Curation par experts seulement (basé sur publis)
  - Attention: toutes les protéines connues n'y sont pas!
    1. Evidence at protein level 72765 13.7%
    2. Evidence at transcript level 69863 13.1%
    3. Inferred from homology 373177 70.1%
    4. Predicted 14474 2.7%
    5. Uncertain 1867 0.4%

# Banques protéiques

- TrEMBL (UniProtKB/TrEMBL):
  - banque protéique produite automatiquement par traduction banque EMBL. 2011: 17.000.000 entrées
- Uniprot=Swissprot+TrEMBL
- Dizaines de Banques spécialisées
  - Cazy (Carbohydrate Active Enzymes)
  - ABC transporters
  - etc.

# NCBI Global query (recherche multi-bases)

Search NCBI databases	
<input type="text" value="beta globin"/>	<input type="button" value="Search"/>
<b>Literature</b>	
<a href="#">8555</a>	PubMed : scientific & medical abstracts/ citations
<a href="#">15605</a>	PubMed Central : full-text journal articles
<a href="#">12</a>	NLM Catalog : books, journals and more in the NLM Collections
<a href="#">4</a>	MeSH : ontology used for PubMed indexing
<a href="#">207</a>	Books : books and reports
<a href="#">28</a>	Site Search : NCBI web and FTP site index
<b>Health</b>	
<a href="#">22</a>	PubMed Health : clinical effectiveness, disease and drug reports
<a href="#">12</a>	MedGen : medical genetics literature and links
<a href="#">15</a>	GTR : genetic testing registry
<a href="#">124</a>	dbGaP : genotype/phenotype interaction studies
<a href="#">124</a>	ClinVar : human variations of clinical significance
<a href="#">122</a>	OMIM : online mendelian inheritance in man
<a href="#">9</a>	OMIA : online mendelian inheritance in animals
<b>Organisms</b>	
<a href="#">0</a>	Taxonomy : taxonomic classification and nomenclature catalog
<b>Nucleotide Sequences</b>	
<a href="#">2182</a>	Nucleotide : DNA and RNA sequences
<a href="#">3</a>	GS5 : genome survey sequences
<a href="#">2042</a>	EST : expressed sequence tag sequences
<a href="#">12</a>	SRA : high-throughput DNA and RNA sequence read archive
<a href="#">29</a>	PopSet : sequence sets from phylogenetic and population studies
<a href="#">124</a>	Probe : sequence-based probes and primers
<b>Genomes</b>	
<a href="#">740</a>	Genome Browser : genome browser interface
<a href="#">100</a>	BLAST : basic local alignment search tool

# Ensembl ([www.ensembl.org](http://www.ensembl.org))

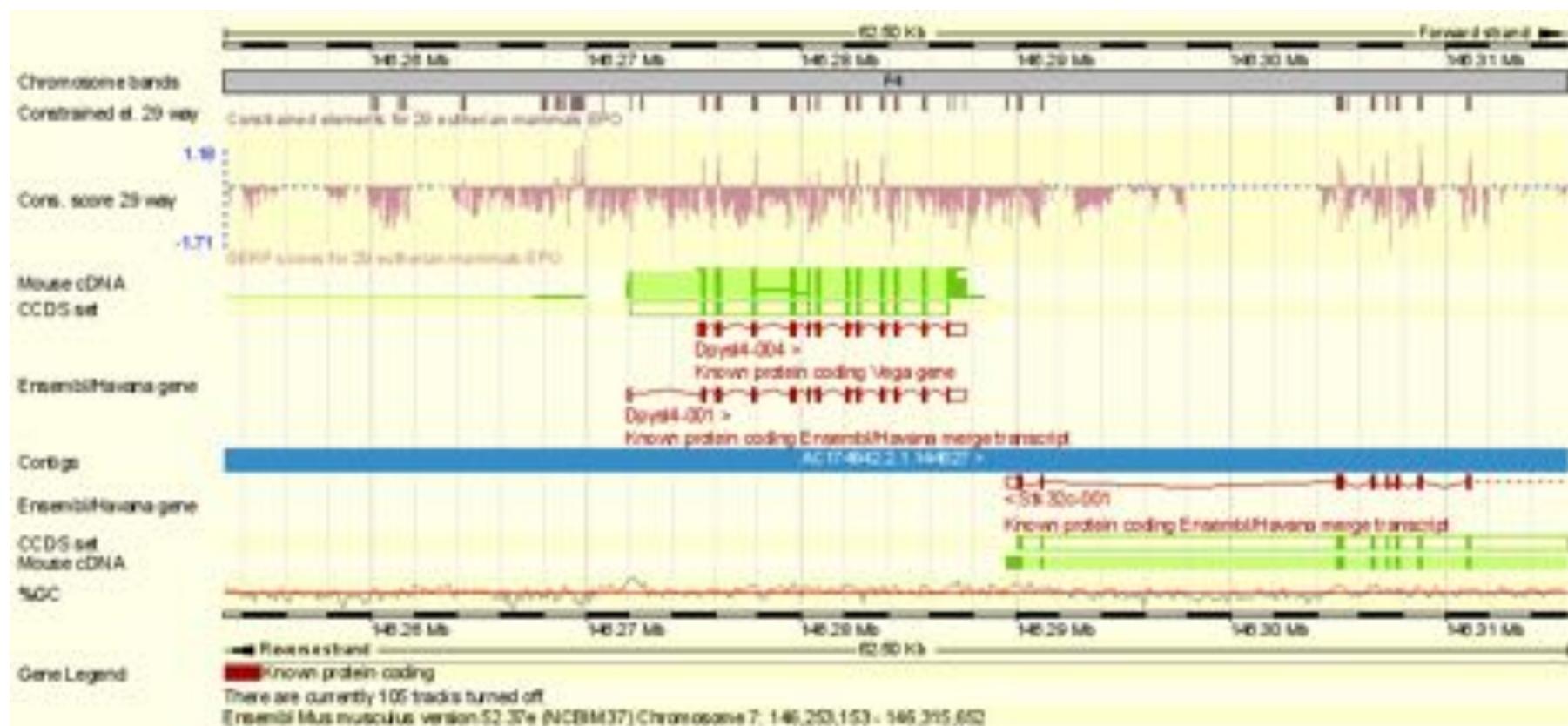
- Plusieurs banques en une:
  - Génomes assemblés
  - Peptides confirmés
  - Transcrits confirmés
  - peptides prédicts
  - Transcrits prédicts
- + un genome browser
- Méthode de prédition (système Genewise): Utilisent programme de prédition Genscan, puis Blast contre: protéines, mRNA, EST
- Génome humain version Sep 2008 (NCBI 36) : Confirmed protein-coding genes: 21649; RNA genes: 4810; Predicted genes (Genscan): 49796; base pairs:  $3,25 \cdot 10^9$ .

Species - Ensembl v24			
Human	prf	NCBI 36	Sep 04
Mouse		NCBIm33	Sep 04
Zebrafish		WTSG.Zv4	Sep 04
Rat		RSGC.R.1	Sep 04
Chicken		WASHUC1	Sep 04
Mosquito		MOZ.2	Apr 04
Fugu		Fugu v2.0	May 04
Frutilly		ICGP 3.1	Jul 03
Chimp		ChIMP1	May 04
Monkeybus		Amel1.1	Sep 04
Tetraodon		TETRAODON7	Sep 04
Dog	prf	BROAD01	
<i>C. elegans</i>		WS_115	Apr 04
<i>C. briggsae</i>		cb05.10cd	Jul 03

# Ensembl Species (2012 - partiel)

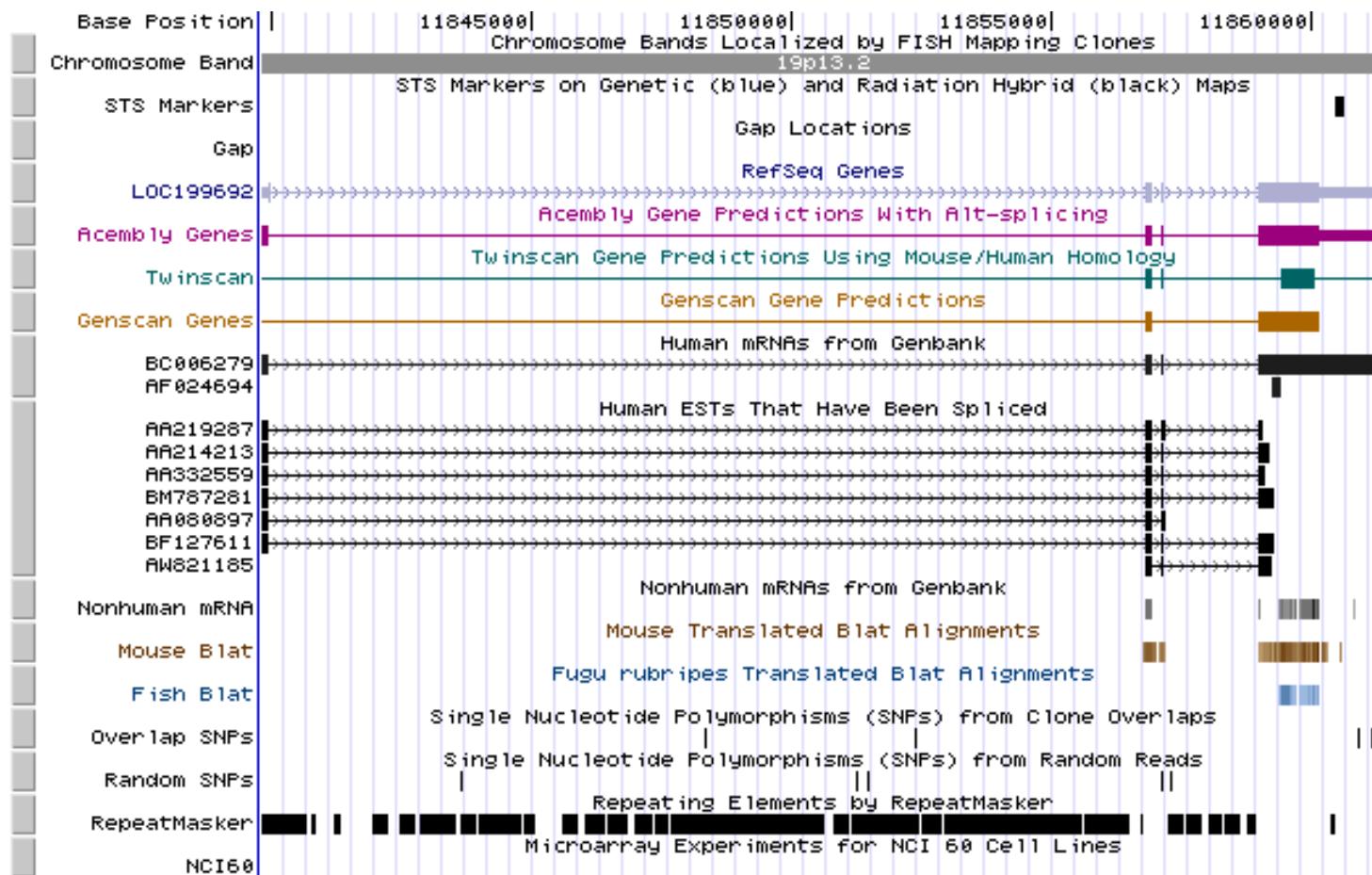
 Alpaca <i>Vicugna pacos</i> v10Pac1	 Gibbon <i>Hominoidea leucogenys</i> HleuG1	 Platypus <i>Oscinopeltidae ornithurus</i> OANAS
 Axole lizard <i>Axolotl</i> AnoCar3.0	 Gorilla <i>Gorilla gorilla</i> gorGor3_1	 Rabbit <i>Oryctolagus cuniculus</i> oryCun2
 Armadillo <i>Taxidea taxus</i> taxtax2	 Guinea Pig <i>Cavia porcellus</i> cavPor3	 Reef <i>Corallium reefense</i> corCor4
 Baboon (древн.-азиатский вид) <i>Papio hamadryas</i> Pham	 Hedgehog <i>Echinaceus europaeus</i> HEDGEHOG	 Sacharomyces cerevisiae <i>Saccharomyces cerevisiae</i> EF4
 Budgerigar (древн.-азиатский вид) <i>Melanisticus undulatus</i> MelUnd3.0	 Horse <i>Equus caballus</i> EquCab2	 Sheep (древн.-азиатский вид) <i>Ovis aries</i> ovsAri1
 Bushbaby <i>Otolemur garnettii</i> OtoGar3	 Hamster <i>Homomys hamster</i> HOMH3T	 Shrew <i>Toromys arenarius</i> COMMON_SHREW1
 Ciona intestinalis <i>Ciona intestinalis</i> CIOIn1	 Hyrax <i>Procavia capensis</i> proCap1	 Sloth <i>Choloepus hoffmanni</i> cholHoff1
 Ciona savignyi <i>Ciona savignyi</i> CIOsav2.0	 Kangaroo rat <i>Dipodomys deserti</i> dipdesert1	 Spotted Gar (древн.-азиатский вид) <i>Lepisosteus osseus</i> LepOsse1
 Caenorhabditis elegans <i>Caenorhabditis elegans</i> WBcel270	 Lamprey <i>Petromyzon marinus</i> Pmarina_7.0	 Squirrel <i>Sciurus vulgaris</i> SciVulgar1
 Cat (древн.-азиатский вид) <i>Felis catus</i> CAT	 Lesser hedgehog tenrec <i>Echinops telfairi</i> TEFREC	 Squirrel monkey (древн.-азиатский вид) <i>Saimiri boliviensis</i> Saimbol1
 Chicken (древн.-азиатский вид) <i>Gallus gallus</i> WASHUC2	 Macaque <i>Macaca mulatta</i> Mmul_1	

# Ensembl: « chromosome view »



# Genome browser UCSC

- <http://www.genome.ucsc.edu/>

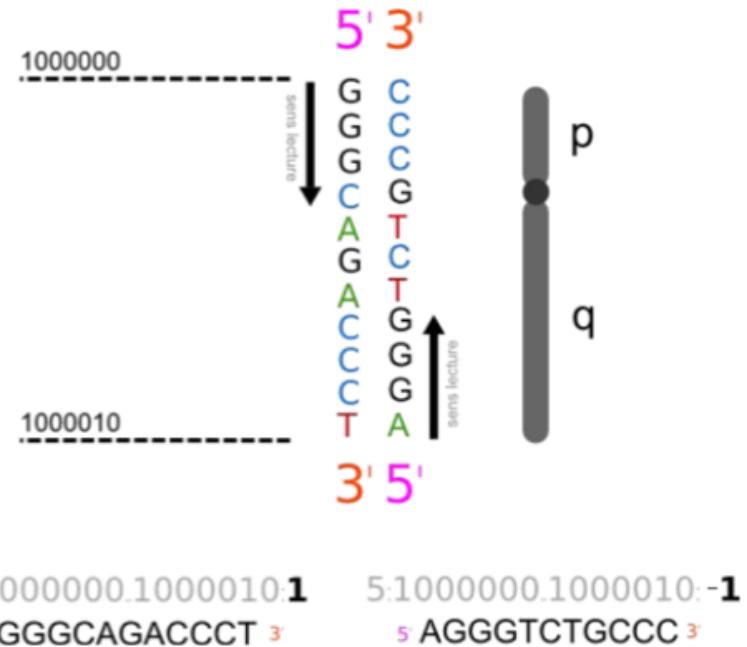


# Accéder à une région génomique avec des API

## API de la banque Ensembl:

<http://rest.ensembl.org/sequence/region/human/7:117465784..117715971:-1>

Attention aux versions  
d'assemblage du  
génome (Hg19, Hg38..)



# Mini TD

- Rendez-vous sur le site [Ensembl](#)
- Cherchez dans le génome humain le gène RPLP0
- Visualisez la région chromosomique, notez les coordonnées.
- Taille du gène? Combien d'introns, exons?
- Récupérer la séquence du mRNA
- Cherchez RPLP0 humain sur le site du [NCBI](#)
- Cherchez RPLP0 humain sur le site [Uniprot](#)
- Avec la [commande wget](#) (curl si vous êtes sur un Mac) et l'API Ensembl, récupérez la séquence du gène RPLP0 humain à partir de ses coordonnées.  
(ajouter options: « ?content-type=text/x-fasta » à l'URL)

# Région chromosomique et annotation transcrits



# Un transcript (= mRNA)

[Ensembl](#) BLAST/BLAT | BioMart | Tools | Downloads | Help & Documentation | Buy | Mirrors

Human (GRCh37) ▾ Location: 12:120,634,468-120,639,014 Gene: RPLP0 Transcript: RPLP0-001

Search all species

Transcript-based displays

- Transcript summary
- Supporting evidence (73)
- Sequence
  - Exons (11)
  - (DNA)
  - Protein
- External References
  - General identifiers (40)
  - Disease pathways (53)
- Ontology
  - Ontology graph (21)
  - Ontology table (25)
- Genetic Variation
  - Variation table
  - Variation image
  - Population comparison
  - Comparison map
- Protein Information
  - Protein summary
  - Domains & features (5)
  - Variations (91)
- External data
  - Protein annotation
- History
  - Transcript history
  - Protein history

Configure this page

Add your data

Export data

Relaunch this page

Share this page

## Transcript: RPLP0-001 (ENST00000382514)

Description: ribosomal protein, large, P0 (Source NCBI: [Symbiosis](#) 18(37))

Location: Chromosome 12:120,634,468-120,639,014 reverse strand

Gene: This transcript is a product of gene [ENSG000000000157](#)

This gene has 27 transcripts (splice variants). [Show transcript table](#)

### Transcript summary

Statistics: Exons: 6 Coding exons: 7 Transcript length: 1,210 bps Translation length: 317 residues

CCDS: This transcript is a member of the Human CCDS ref: [CCDS00100](#)

Ensembl version: ENSEMBL 73 (GRCh37) 14.4

Type: Known protein coding

Prediction Method: Transcript where the Ensembl genebuild transcript and the [UCSC](#) manual annotation have the same sequence, for every base pair. See [gtf2gtf](#)

Alternative transcripts: This transcript corresponds to the following database identifier(s):

Transcript having exact match between ENSEMBL and HAVABA: [GTHUMT00000402458](#) (version 2)

Ensembl release 73 - September 2013 © EMBL / EBI

About Ensembl | Privacy Policy | Contact Us

Human (GRCh37) · Gene to protein site

# Les variations du génome

# En quoi nos génomes diffèrent-ils?

- Les individus d'espèces différentes ont des génomes différents par la taille, l'ordre et la nature des informations qu'ils contiennent.
- On considère souvent que 2 individus de la même espèce possèdent le « même » génome.
- En fait, le génome de chaque individu est unique. Chez l'homme le génome diffère de 0.1% entre 2 personnes non apparentées

# Les variations du génome dans une population

- **Très importantes médicalement**
  - Pharmacogénomique: comment chaque patient répond aux drogues
  - Marqueurs de susceptibilité aux maladies
- **Polymorphismes dans le génome humain**
  - Insertions, délétions, duplications, réarrangements
  - Microsatellites etc..
  - Single Nucleotide Polymorphism (SNP)
    - Les plus fréquents

# Les SNP: Single Nucleotide Polymorphism

Les polymorphismes les plus fréquents dans la population

SNPs or SNPs =

sites of variation in the genome  
(spelling mistakes)

Karen	AGCTTGAC TCCATGATGATT
Debo	AGCTTGAC <b>GCC</b> ATGATGATT
Jose	AGCTTGAC <b>TCC</b> C TGATGATT
Thomas	AGCTTGAC <b>GGCC</b> TGATGATT
Anupriya	AGCTTGAC TCCATGATGATT
Robert	AGCTTGAC <b>GCC</b> ATGATGATT
Michelle	AGCTTGAC <b>TCC</b> C TGATGATT
Zhijun	AGCTTGAC <b>GGCC</b> TGATGATT

Dans le génome humain: **1 SNP chaque 500 bp** → ~6 million

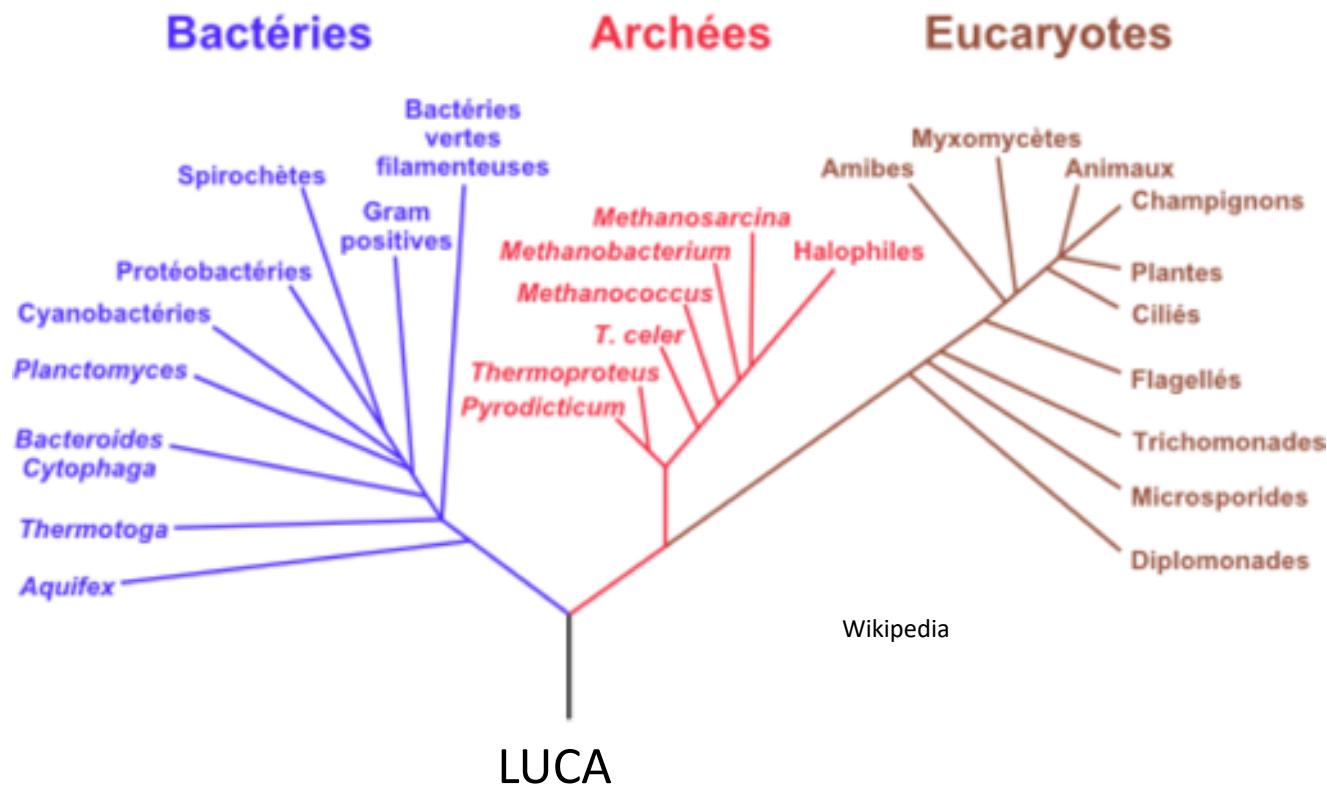
L'ADN de deux individus distincts est à 99,5% identique

[cpmc.coriell.org/Sections/Medical/GenelInterac.](http://cpmc.coriell.org/Sections/Medical/GenelInterac.)

# La ressemblance entre génomes

- Homme A / homme B
  - 99,9% identique (0,1% différence)
- Homme/chimpanzé
  - Codant: 98,5% identique
  - Non codant: ~96% identique
  - 90Mb d'insertions/délétions et 35 millions de différences ponctuelles
- Homme/souris
  - Codant: 90% identique
  - Non codant: la majorité est sans identité apparente, mais on trouve quand même de nombreux segments semblables (« conservés »)
- Homme/poulet
  - Codant: 80% identique
- Homme/poisson
  - Codant: 70% identique

# Rappel: l'arbre du vivant



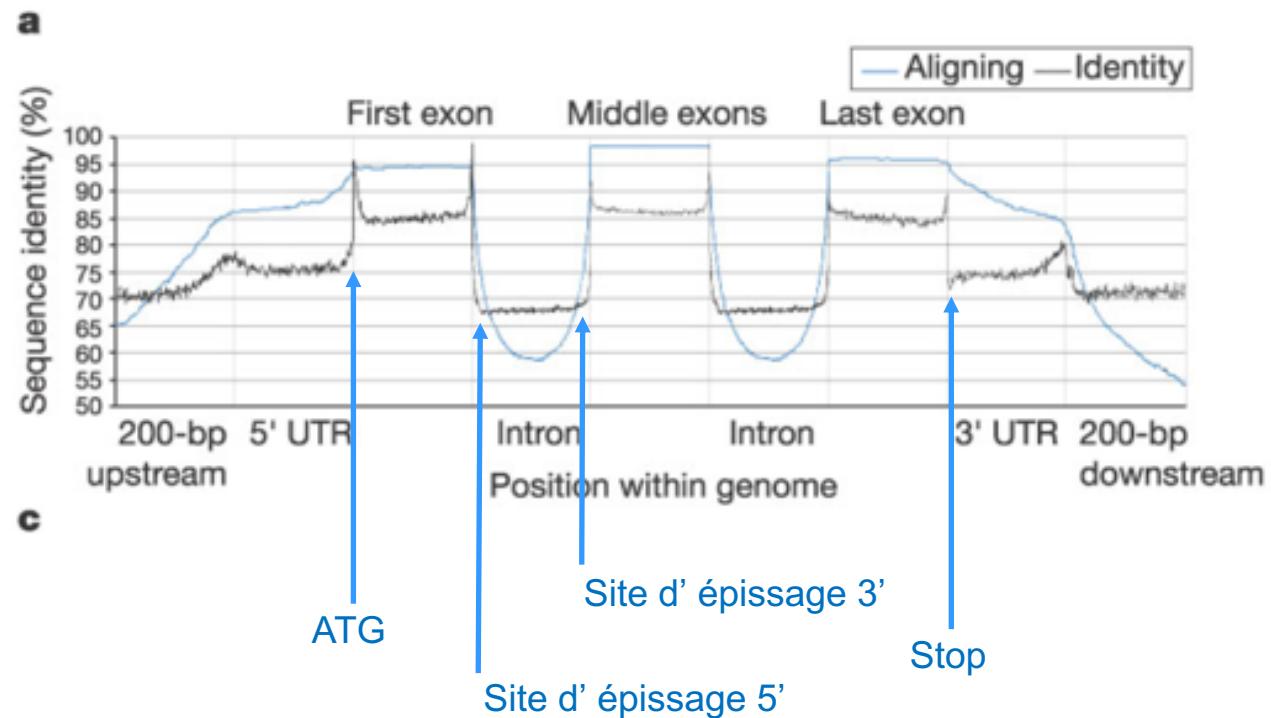
# La conservation des séquences codantes ou non

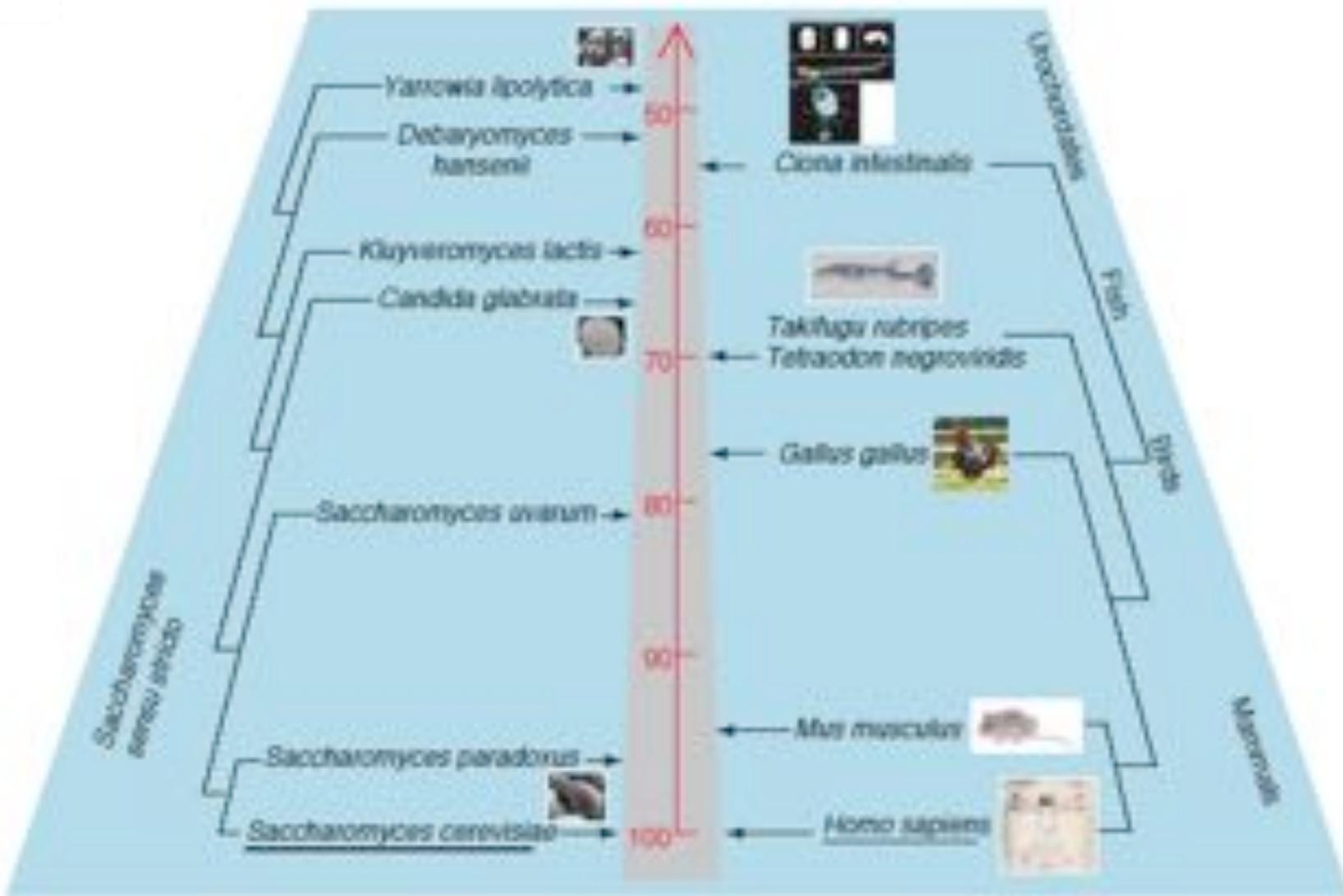
Figure de l'article de 2002 sur le génome de la souris.

“a wonderful visual guide to the most important features of mammalian genes” (Chris Ponting)

Identité homme-souris pour tous les couples de gènes orthologues

Divergence : 90M années





From B. Dujon, Trends in Genetics, 2006  
Diatomées: Chris Bowler

# Comparer les séquences

# RPLP0: un gène de protéine ribosomale universellement conservé.

mRNA humain de RPLP0 (exons seulement)

- >ENST00000392514 cdna:KNOWN\_protein\_coding  
GTCTGACGGCGATGGCGCAGCCAATAGACAGGAGCGCTATCCGCGGTTCTGATTGGCT  
ACTTTGTTCGCATTATAAAAGGCACGCGGGCGGAGGCCCTCTCGCCAGGCGTCC  
TCGTGGAAGTGACATCGTCTTAAACCCTGCCTGGCAATCCCTGACGCACGCCGTGATG  
CCCAGGGAAGACAGGGCGACCTGGAAGTCCAACACTTCCTTAAGATCATCCAACATTG  
GATGATTATCCGAAATGTTCAATTGTGGGAGCAGACAATGTGGGCTCCAAGCAGATGCAG  
CAGATCCGCATGTCCCTCGCGGGAAAGGCTGTGGTGTGATGGCAAGAACACCATGATG  
CGCAAGGCCATCCGAGGGCACCTGGAAAACAACCCAGCTCTGGAGAAACTGCTGCCTCAT  
ATCCGGGGAATGTGGGCTTGTGTTACCAAGGAGGACCTCACTGAGATCAGGGACATG  
TTGCTGGCCAATAAGGTGCCAGCTGCTGCCGTGCTGGGCCATTGCCCATGTGAAGTC  
ACTGTGCCAGCCCAGAACACTGGTCTCGGGCCCGAGAAGACCTCCTTTCCAGGCTTA  
GGTATCACCAACTAAAATCTCCAGGGCACCATTGAAATCCTGAGTGATGTGCAGCTGATC  
AAGACTGGAGACAAAGTGGAGGCCAGCGAAGCCACGCTGCTGAACATGCTAACATCTCC  
CCCTCTCCTTGGCTGGTCATCCAGCAGGTGTTGACAATGGCAGCATCTACAACCCT  
GAAGTGCTTGTATATCACAGAGGAAACTCTGCATTCTCGCTTGGAGGGTGTCCGCAAT  
GTTGCCAGTGTCTGTGAGATTGGCTACCCAACTGTTGCATCAGTACCCATTCTATC  
ATCAACGGGTACAAACGAGTCCTGGCCTTGTCTGTGGAGACGGATTACACCTCCCCACTT  
GCTGAAAAGGTCAAGGCCTTCTGGCTGATCCATCTGCCTTGTGGCTGCTGCCCTGTG  
GCTGCTGCCACCACAGCTGCTCCTGCTGCTGCAGCCCCAGCTAACGGTTGAAGCCAAG  
GAAGAGTCGGAGGAGTCGGACGAGGATATGGGATTGGCTCTTGACTAACCAAAA  
AGCAACCAACTTAGCCAGTTTATTGAAAACAAGGAAATAAAGGCTACTTCTTAAA  
AAGTCTCTGGACTCTTAA

# Alignement avec mRNA RPLP0 de *Pan troglodytes*

GENE ID: 462901 RPLP0   ribosomal protein, large, P0 [Pan troglodytes]			
Score = 2141 bits (2386), Expect = 0.0 Identities = 1211/1218 (99%), Gaps = 1/1218 (0%) Strand=Plus/Plus			
Query	1	GTCCTGAGGGGGATGGGCGCAACAAATGACAGGGAGGCTA200GCGGGTTCTGATTCGT	60
Subject	1	GTGGGATGGGGGGATGGGCGCAACAAATGACAGGGAGGCTA200GCGGGTTCTGATTCGT	60
Query	61	ACTTTGGTGGATTAATAAAAGTCAGGCGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG	120
Subject	61	ACTTTGGTGGATTAATAAAAGTCAGGCGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG	119
Query	121	TCTTGGAAGTCACATGCTTTAAACCCCTGCATGGCAATGCCCTGAGCGACCGCGCTGAT	180
Subject	120	TCTTGGAAGTCACATGCTTTAAACCCCTGCATGGCAATGCCCTGAGCGACCGCGCTGAT	179
Query	181	CCCCGGGAAAGCAGGG	240
Subject	180	CCCCGGGAAAGCAGGG	239
Query	241	GATGATTAACTGGAAATTTTCACTTGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG	300
Subject	240	GATGATTAACTGGAAATTTTCACTTGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG	299
Query	301	CAGATCGGCGATGTTGG	360
Subject	300	CAGATCGGCGATGTTGG	359
Query	361	GGCAAGGGCGATCGGAGGG	420
Subject	360	GGCAAGGGCGATCGGAGGG	419
Query	421	AATGGGGGGAAATGTGGGCTTTGTGGTCAACGAGGGGGCTCACTGGGATCAGGGGACATG	480
Subject	420	AATGGGGGGAAATGTGGGCTTTGTGGTCAACGAGGGGGCTCACTGGGATCAGGGGACATG	479
Query	481	TTTCTGGGCAATTAATGTTGGCGCTTGCTGGCTGATTCGGGGGGGGGGGGGGGGGG	540
Subject	480	TTTCTGGGCAATTAATGTTGGCGCTTGCTGGCTGATTCGGGGGGGGGGGGGGGGGG	539
Query	541	ACTATGGCGAGGG	600
Subject	540	ACTATGGCGAGGG	599
Query	601	GGTATCACCACTAAATCTCGAAGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG	660
Subject	600	GGTATCACCACTAAATCTCGAAGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG	659

# Alignement avec mRNA RPLP0 de *Danio rerio*



> ref NM_131602.3   D Danio rerio ribosomal protein, large, P0 (rplp0), mRNA  gb BC042854.1   UGM Danio rerio ribosomal protein, large, P0, mRNA (cDNA clone MDC177791  IMAGE:7001273), complete cds Length=1105	
GENE ID: AB101 rplp0 : ribosomal protein, large, P0 (Danio rerio) (Over 10 Published links)	
Score = 388 bits (984), Expect = 0.0 Identities = 801/1000 (80%), Gaps = 10/1000 (1%) Strand=Plus/Plus	
Query 137	GTCTTTAACCC--CTGCTTGCAATCCTG-ACGCCACGCCCTGATGCCCGAGGAAACA 193
Subject 91	GTCTTTAACCCGCTCTTCAACGAACTCTGAAAGCACTGCAAGGATGCCCGAGGAAACA 90
Query 134	GCGGGAACGTTGAAAGTCCTAACACTCTCTTAAGATCATGCAACTATTGGATGATTATCGA 253
Subject 91	GCGGCAACGTTGAAAGTCCTAACACTCTCTGAAAATCATGCAACTGCTGGATGACTACCCA 189
Query 254	AAAGTTTCATTGTTGGGAGCAGAACATGTTGGGCTGAAAGCAGATGCGAGCAGATGGCAATG 313
Subject 151	AGTGTTCATGTTGGGCGAGCAGAACATGTTGGGCTGAAAGCAGATGCGAGCAGATGGCAATG 210
Query 314	CCCTTGCGGGAGGGCTGTGGTGTGATGGGCAAGAACACCATGATGGGCAAGGCATCC 373
Subject 211	CCCTTGCGGGCAAGGGCTGTGGTGTGATGGGCAAGAACACCATGATGGGCAAGGCATCC 270
Query 274	GAAGGCACCTGAAAGAACACCCAGCTCTGGAGAAACTGCTGCTCAATGCGGGGAAATG 493
Subject 171	GTGGGCACCTGAAAGAACACCCAGCTCTGGAGAGGGCTGCTTCCGACATGCGGGGAAACG 390
Query 434	TGGGGCTTGTCTTCAACGAAAGGAGGCTCACTGAGATGAGGACATGTTGCTGGGCAATA 493
Subject 331	TGGGGCTTGTCTTCAACGAAAGGAGGCTGAGATGAGGCTGGGAGGCTGCTGCTGGGCAATA 390
Query 494	AAATGCCACCTGATGCCCTTGCTGCTGCTGCCATGTTGAGCTACTGCTGCCAGCCC 583
Subject 391	AAATGCCACCTGATGCCCTTGCTGCTGCTGCTGCCATGTTGAGCTGACGCTGCCAGCCC 460
Query 554	AGAACACTGCTGCTGGGCGAGAAGAGGCTCTTTCAGGCTTTAGGTATCACCACTA 613
Subject 451	AGAACACCGGGCTGGGCTGAGAAGAGGCTCTTTCAGGCTTTGGGAAATCACCACTA 510
Query 614	AAATCTGAGGGGCAACATTGAAATCTGAGTGTGCAAGCTGATCAAGAGCTGGAGGACA 673
Subject 511	AAATCTGAGGGGCAACATTGAAATCTGAGTGTGCAAGCTGATCAAGAGCTGGAGGACA 510

# Alignement avec mRNA RPLP0 de *Drosophila melanogaster*



# Alignement avec mRNA RPLP0 de *Arabidopsis thaliana*



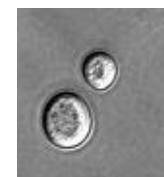
>RefSeq NM\_129359.2 | UEGM | **Arabidopsis thaliana** 60S acidic ribosomal protein P0-1 (AT2G40010) mRNA, complete cds  
Length=1163

GENE ID: NM\_129359 AT2G40010 | 60S acidic ribosomal protein P0-1 (Arabidopsis thaliana) [10 or fewer Published Links]

Score = 66.2 bits (72), Expect = 6e-08  
Identities = 261/402 (65%), Gaps = 14/402 (3%)  
Strand=Plus/Plus

Query	Subject	Score
409	CTGCCTGCTCATATCCGGGGAAATTTGGCCTTGTTCACCAAGGAGGACTCTGAG	468
341	CTGCCTGCTCTCTTCAGGGGAAATTTGGTGACTTCACCAAGGAGGACTCTGAG	350
469	AATGAGGGAC-AAGTTTGTGGGCGATAAAGGTTGGCAGCTGGCTGGGCGGGGGATTGC	527
321	GTCAGTGAAGAGTTGTAACTAC-AAGGTTGGAGCTGCTGATGTTGAGGTTAGTGC	359
628	CCGATGGAAGTCACTTGCGAGCGAG--AACACTGTTCTGGGGGGGGAGAGAGACCTGC	555
360	TCCAATTAGTGTGTTGTCGAA--CGGGGAGAGCTGTTGACGTTTCACGAGACCTGC	417
586	TTTTTCCAG--GCTTAGTTATCAACCAATAAATCTCGAGGGCACCCATTGAAATCTGA	643
418	TCCTTCCAGGTGCTTAA--CTTCCAAACCAAATCAGAAAGGTTGAGATCATAA	475
644	GTGATGTGAGGTGATCGAGAAGCTGGAGACAAAAGTGGGAGGCCAGGCGACGCGCTGCTGA	703
476	CCCGTGTGGAGCTCATCGAAASGGCGACAAAGTGGGTTGATCCGAGGCTGCGCTTCTG	536
704	ACATGCTCAACATCTGGGGCTCTGCTGGCTGGCTGATOCAGCAGGTGTTGACAATG	743
536	CCAAAGCTTGGAAATCGGGCTTCTGATGATGCTGTTGAGTCACTGAGGATAATG	695
764	G--CAAGCTTACACCGTGAAGTGTGATATCACAGAGGA	859
896	GTCAG--GTTTAACCGTGAAGTGTAACTCACTGAGGA	635

# Alignement avec mRNA RPLP0 de *Schizosaccharomyces pombe*



# RPLP0: la protéine

>sp|P05388|RLA0\_HUMAN 60S acidic ribosomal protein P0 OS=Homo sapiens GN=RPLP0 PE=1 SV=1  
MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKQMQQIRMSLRGKAVVLMGKNTM  
MRKAIRGHLENNPALEKLLPHIRGNVGTVFTKEDLTEIRDMLANKVPAARAGAIAPCE  
VTVP AQNTGLGPEKTSFFQALGITTKISRG TIEILSDVQLIKTGDKG VASEATLLNMLNI  
SPFSFGLVIQQVFDNGSIYNPEVLDITEETLHSRFLEGVRNVASVCLQIGYPTVASVPHS  
IINGYKRVLALS VETD YTFPLAEKVK AFLADPSAFVAA APVAA ATTAA APAAA A PAKVEA  
KEESEESDEDMGFGLFD



>sp|P19889|1|SLAO\_DROSOPHILA  
RecName: Full=60S acidic ribosomal protein P0; AltName: Full=apurinic endonuclease; AltName: Full=DNA-(apurinic or apyrimidinic site) lyase  
Length=317

Score = 417 bits (1072), Expect = 2e-146, Method: Compositional matrix adjust.  
Identities = 210/317 (66%), Positives = 205/317 (65%), Gaps = 0/317 (0%)

Query 1	MPPREDRATKNSHNTFLKTIQQLLDGYPRCFCFIVGADMSVGSKQHQIRNMLPCKWAVYVLMGRKHTM	60
Subject 1	MVRERKKAAMKKAQYTFIIVVVELFDDEFPRCFCFIVGADMSVGSKQHQIRNMLPCKWAVYVLMGRKHTM	60
Query 41	MRRKAIRGNLLENNPALEKLLPPIRGMNVGVFVFTKEDLTKEIRDMILLANEVPAARAGAIAPCE	120
Subject 41	MRRKAIRGNLLENNP LERKLLPPI+CNVGVVFTR DL E+RD LL +EV A AP GAIAAP	120
Query 121	VTVPAQNTGLGPENTEFFQALGTTKISGTTTIEILSDFVQLIETGDKVNGASKEATLLNLNT	180
Subject 121	VVTPAQNTGLGPENTEFFQAL I TKG+GTIEI++EV ++K GOKVGAKEATLLNLNT	180
Query 181	SPPFEGCLVTDQQVYDQGSIYTMPEVLDITEXTTMRFLQCVPMVAVYVCLQIDGYPTVAEVPMH	240
Subject 181	SPPFEGCLIVNDQVYDQGSIYTMPEVLDIINPVEDLAAKFOQCVPMVAVYVCLQIDGYPTIAEVPMH	240
Query 241	IINGEKRVIALSVEVDYTFPLAKKNAFLADPRAVYAAAPVRAATTAAAPAAAAPVAYEA	300
Subject 241	IANGFIOTLALIAATTEVEPKEATTIKEYEINDPKEFAAAFAAAAAPANGGATEKKEEAEKP	300
Query 301	KEEESEEESDDEIDNGFGLFID 317	
Subject 301	KEEESEEESDDEIDNGFGLFID 317	



>sp|P67491|18LAD03 ARATH **G** DecName: Full=60S acidic ribosomal protein P0-3  
Length=323

GENE\_ID: 320235 ATCGG11280 : 60S acidic ribosomal protein P0-3  
(Arabidopsis thaliana) (10 or fewer PAMAC links)

Score = 311 bits (798), Expect = 1e-104, Method: Compositional matrix adjust.  
Identities = 168/321 (52%), Positives = 220/321 (69%), Gaps = 4/321 (1%)

Query 1	MPEGEDGATMMKMYFLKTIQILLLDQYVWFCFTVQADAFVVAQHMQGTRHEMLRQKQVVLMSKMTK	60
	M + + A K Y E+ QL+D+Y + + V ADRYVGS Q+Q IR LRG +VVLMSKMTK	
Subject 1	MVKATKAERKIAVDTKLQQLIDEYTQILVVAADMVQSTQLQMKRKGSLAGGQVVLMGKMTK	60
Query 41	MKKAIRGHLME--MPALEKLLPMPHRGDNVQFVFTKEDLTKEIDMLANKVPAAAARAGAAIAP	118
	M++++R H DK H A+ LLP ++GRVG +FTK DL E+ + + EV A AR G +AP	
Subject 41	MKQASVRIHRSSENQNTAIINLLPLLGQMNVLIFTKQDLKPVWEEVAKVYVQGAPAKVQELVAP	120
Query 119	CETVVRQAQNTGLGPDKTISFFQALGITTNIKISRGTKIEILSUVQQLIKTGOWVGAKEATLML	178
	4V V NTGL P +TEFFQ L I TKI++GT+K++ V+LIK GOWV+SEA LE L	
Subject 121	IDVVVQPGFTGLOPSPQTRFFQVLMNPTKINNGTIVKIIITWVELIKQGOWVGGSEAAALLAKL	180
Query 179	HSIPPFETGLVVIQGVFTDNGSITMPEVLDITETLNSRPFLEGTVANVAVSVCLQDGYPTVAEVF	238
	I DPF+GLV+Q V+DNQ++4PEVLD+TE+ L +F Q+ V S+ L + YPT+A+ P	
Subject 181	GIRPFETGLVVIQGVFTDNGSVEPEVLOLTDQLVETFASCIEMNTELALAVSYPTLAAAAP	240
Query 239	HSIIMGYTQVIALSVEVDYTFPLAENKVAFLADPDAIFVAAAPVAAATTAAAPAAALPAAVY	298
	K DK YK L+V V DQVFF AAEVK F+ DPS TV AA +A +A A A	
Subject 241	HSIFINAYKHALAIVATDVTFPQAEKVVAFLDPSKFTVAAATTAADAGGGGAAQAAJAK	300
Query 299	EAKKEESEEEDDEIM--GFGLFQ 317	
	+++ E +ED GFGLFQ	
Subject 301	VEEKKKEESEEEDYEGGGFGLFD 321	



Map:074864.1|SILAO\_SCP0

G RecName: Full=60S acidic ribosomal protein P0

Length=312

Gene ID: 2538823 rps0 | 60S acidic ribosomal protein Rps0 (predicted)  
(Schizosaccharomyces pombe 972h-1) (10 or fewer Published links)

Score = 318 bits (816), Expect = 1e-107, Method: Compositional matrix adjust.  
Identities = 170/308 (55%), Positives = 219/308 (71%), Gaps = 3/308 (1%)

Query 10 KENPYPLATIQILLDOYPRCFTVQADNVGSKMHQQISDMSLRSACAVVILMHDNTMGRKAIDQHL 69  
H+ YF H+ L + T T+V DMV S+QM +R LRG A ++NGHTHM-R+A+DG +  
Subject 8 KAQYPERKLRLSPKTYNPKLPPVVMIDAVV185Q-QGIVTVAQQLRSTAKLIMDQHNTMIRGAMGII 67

Query 70 EMDPDALEKLLPHIIGRQW/GFVFTKEDLTKEIICMLLANKVPAAAADAGATAIPCEVTVRAQMTG 129  
+ P LE+LLP +RCMVGFVFT DL E+R+ ++AK + A AR AIAA +V VFA NTG  
Subject 48 NMOPKELERLLPVV/RGMW/GFVFTNAQKLVKEVTTIIASVIAAARPMKIAIPDQVTFVRAQMTG 127

Query 130 LGPDKTSYYFQALGITTIXISPGTIEKILSDVQQLIKTCGDTVCAEATLLKGKISFFFPGLV 169  
+ P KTSTFQALGI TKI+PGTIEK SDV 14 KTG SEATLLKGKISFFF+Q+ \*  
Subject 128 MHPDKTSYYFQALGIFTIXISPGTIEKILSDVQQLIKTCGDTVCAEATLLKGKISFFFPGLV 167

Query 190 QQVYDQGSEIYINPHEVLDITETLMPFLEGVYRIVVAASVYCLQ1GYPVVAJFVWHSIINGTKPVL 249  
++D G++++DE+LD++EE L I + ++ L TPT+ SV ED+EE TK ++  
Subject 188 LTIYDQGNTVTSRKEILDVKEEDLIGHILSAASITTAISLGANTFTILDF3HGSVTVWATVNLV 247

Query 250 ALSVETDYYFPLAIDKWAFLADPSAFVAAAPVAAATRAPAAAAAAPAIVVAKKESSEEDE 309  
A+S+ T+YF E+ KAFLADPSAFV A AA AA A APA A EE EEDDE  
Subject 248 AVSLATEYYFESTEQKWAFLADPSAFVVA---APRAAAAGGEAAAPAAAAAEEEEEDE 304

Query 310 DMGFGGLFD 317

Subject 306 DMGFGGLFD 312

UniProtKB/Swiss-Prot P01787 L10E0\_2HICR SecName: Full=Acidic ribosomal protein PO homolog; AltName: Full=L10E Length=339  
**Gene ID:** P01787 **UniProtKB/Swiss-Prot** | acidic ribosomal protein PO  
*(Thermococcus onnurineus 5A1)* (10 or fewer PubMed links)

Score = 131 bits (353), Expect = 4e-38, Method: Compositional matrix adjust.  
 Identities = 77/257 (30%), Positives = 130/257 (51%), Gaps = 3/257 (1%)

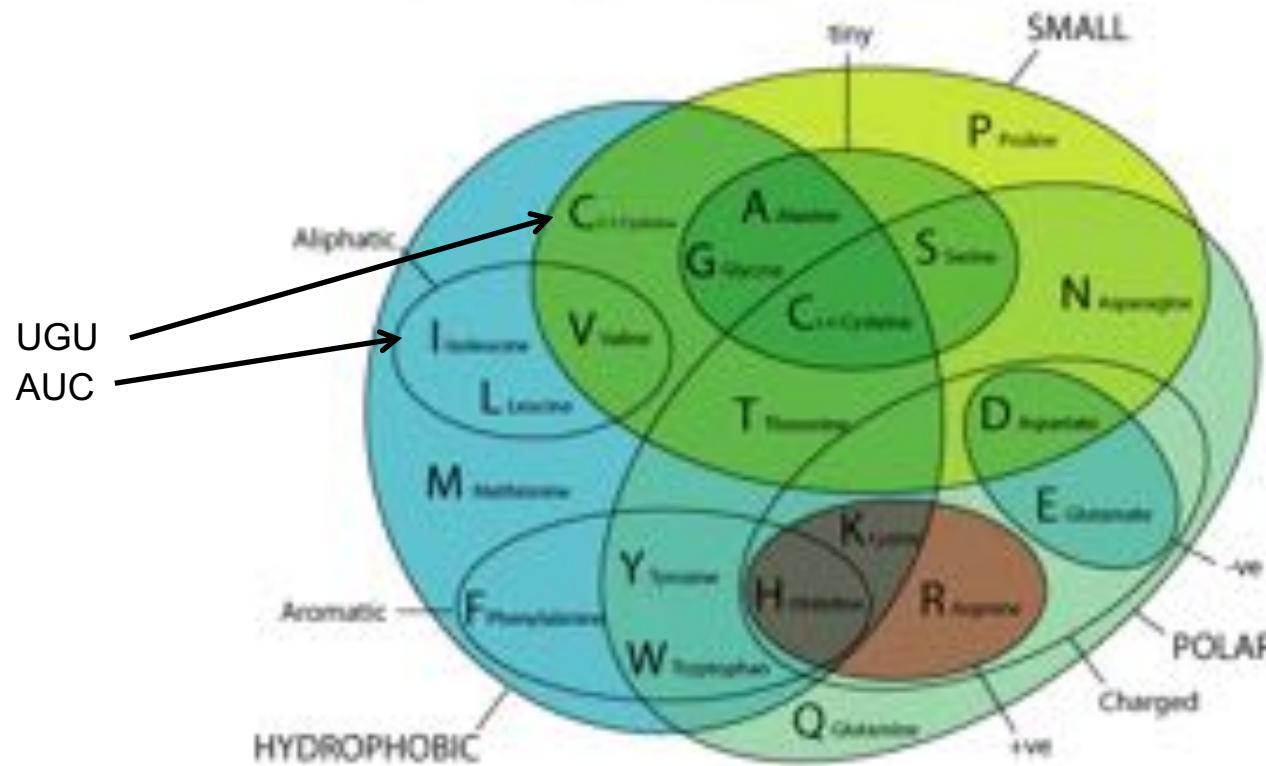
Query	7	ATWKSNTTFLKIIQLLDDTYPRCFTIVGADNPVGSNQMQQTAMNLGRKATVLMGNVTHMRKAIR	66
	A MK	++ ++ YP +U SV + * + +E LRQKA++ + 4ST++ AI+	
Subject	8	ADMKKEEVHELTWITKHYVVIALIVNAVNWAVYFLSNGRECLSGALLVEVNTLILAEIK	64
Query	67	GHLN--WPALEKLLPWRHGMVGFVPTWEDLTIKIRDMLLANKVPAAAARAGAZAPCEVTVP	124
	+ P LEKL+ H+S G + T+ * ++ +L +K SR A+ G P +V +F		
Subject	68	RAAQELGKPELEKLLIDHIQGGAGILATEHSPFWLYKLLERETPAPANGGV/PVFRDVVIP	124
Query	126	AQMTGLOOEEK-TSFFPALGITTWSRGTTIELSDVQLIKTGDNNGASEATLISGLNISRF	183
	A T + P QALGI +I +G + I D ++K G+ + A +IN L I P		
Subject	126	AGPTTSISPGSLVGEHQALGTTTRAKIEKGKVSIQKDVTVLKAGEVITEQLARTIMALGIEL	184
Query	184	EVGVVVIQQYVDMGKTTMVEVDITTEETLNSRLTLEGVYRMAASVCLQIGSTVAVSYWNSIIN	243
	GL + ***G +Y P EVL I EE * * * + + * YPT +* I		
Subject	186	EVGVVLLAAAYEDGIVYTREVIAIDDEKEYIYLQQAIIHAFKLSTNTATFTSQTIKAI:QK	244
Query	244	GYKMYLALSIVETDTTFP	240
	Y +***Y P		
Subject	245	AYLGAGDFVAWEAGYITP	241

# Pourquoi la comparaison de protéines est-elle plus sensible que la comparaison d'ADN?

le code génétique										
	Deuxième lettre									TAC
	U	C	A	G						
U	UUU	Phe	UCU	Ser	UAU	Tyr	UGU	Cys	U	
	UUC	Phe	UCC	Ser	UAC	Tyr	UGC	Cys	C	
	UUA	Leu	UCA	Ser	UAA	Stop	UGA	Stop	A	
	UUG	Leu	UCG	Ser	UAG	Stop	UGG	Trp	G	
C	CUU	Leu	CCU	Pro	CAU	His	CGU	Arg	U	Troisième lettre (côté 3')
	CUC	Leu	CCC	Pro	CAC	His	CGC	Arg	C	
	CUA	Leu	CCA	Pro	CAA	Gln	CGA	Arg	A	
	CUG	Leu	CCG	Pro	CAG	Gln	CGG	Arg	G	
A	AUU	Ile	ACU	Thr	AAU	Asn	AGU	Ser	U	
	AUC	Ile	ACC	Thr	AAC	Asn	AGC	Ser	C	
	AUA	Ile	ACA	Thr	AAA	Lys	AGA	Arg	A	
	AUG	Met	ACG	Thr	AAG	Lys	AGG	Arg	G	
G	GUU	Val	GCU	Ala	GAU	Asp	GGU	Gly	U	
	GUC	Val	GCC	Ala	GAC	Asp	GGC	Gly	C	
	GUA	Val	GCA	Ala	GAA	Glu	GGA	Gly	A	
	GUG	Val	GCG	Ala	GAG	Glu	GGG	Gly	G	
codon d'initiation					codon de terminaison					

## Deuxième raison (plus importante):

### Amino Acid Properties

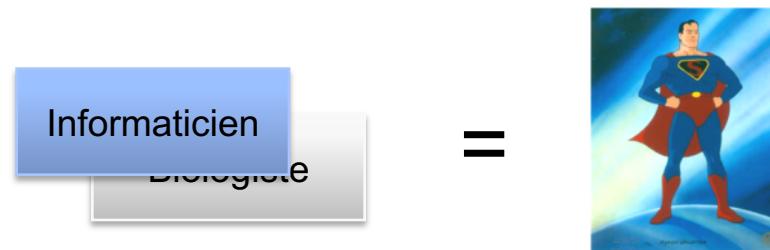
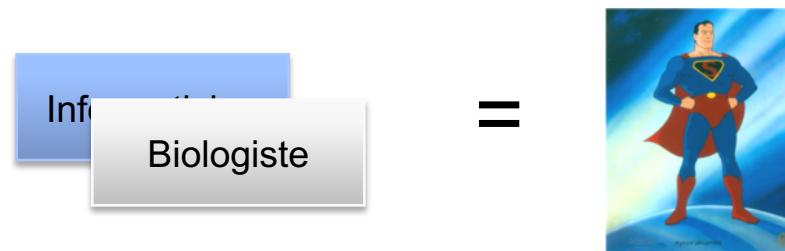
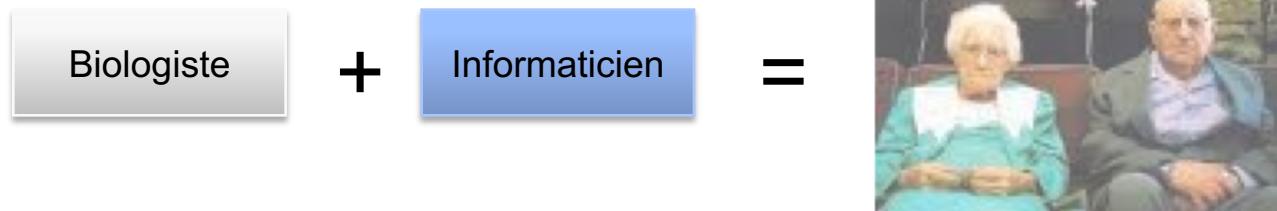


From Livingstone, C. D. and Barton, G. J. (1993),  
"Protein Sequence Alignments: A Strategy for the Hierarchical Analysis of  
Residue Conservation", Comp. Appl. Bio. Sci., 9, 745-756.

-> nécessité de capturer ces propriétés dans un score

# Bioinformatics Wisdom

# Bio-informatique et bioinformatique



# Another version...

