

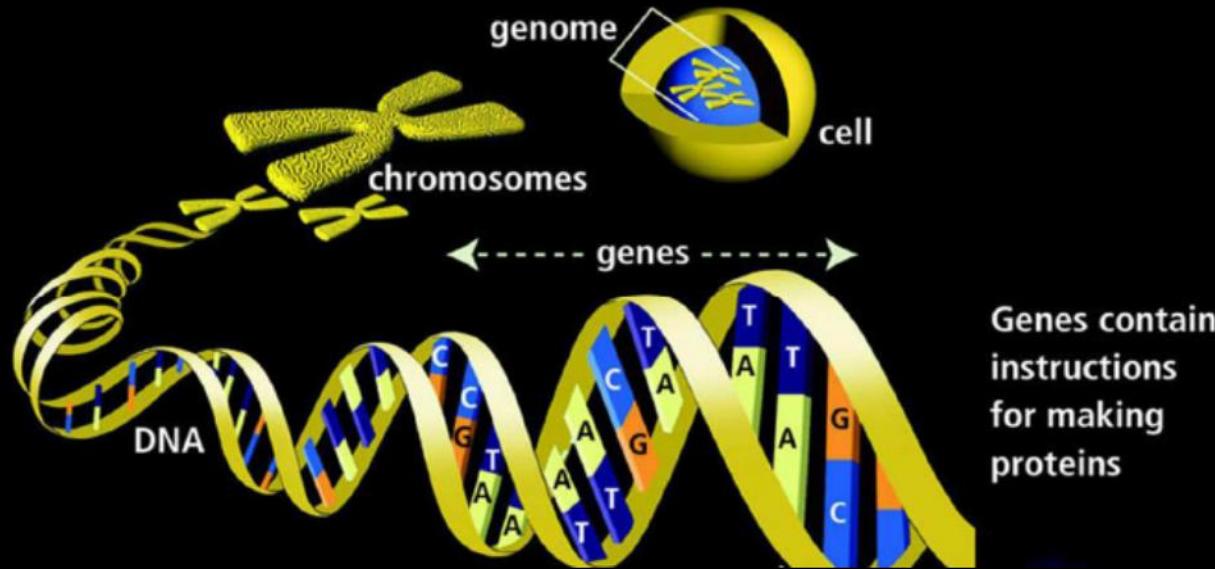
# Introduction to Genomics

Sébastien Aubourg



# Genome

---



Each cell contains a complete copy of the genome,  
distributed along chromosomes  
( $3 \times 10^9$  base pairs - 3 Gb- in human DNA: 6 meters in each cell)

Encodes blueprint for all cellular structures and activities  
and which cells go where...

# Genome

---

1920 : All genes of an organism

Now, the genome is all the DNA in a cell :  
All the DNA on all the chromosomes  
Includes genes, intergenic sequences, repeats

Eukaryotes can have 2-3 genomes

- Nuclear genome
- Mitochondrial genome
- Plastid genome

If not specified, "genome" usually refers to the nuclear genome.

# Genomics

---

Genomics is the study of genome(s), including large chromosomal segments containing many genes.

The initial phase of genomics aims sequence an initial set of entire genome(s).

Functional genomics aims to deduce information about the function of DNA sequences (structural and functional annotations from bioinformatics and high-throughput experimental approaches). Should continue long after the initial genome sequences have been completed...

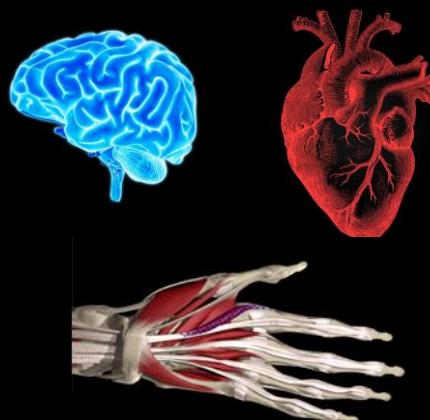
# Different genomes

---

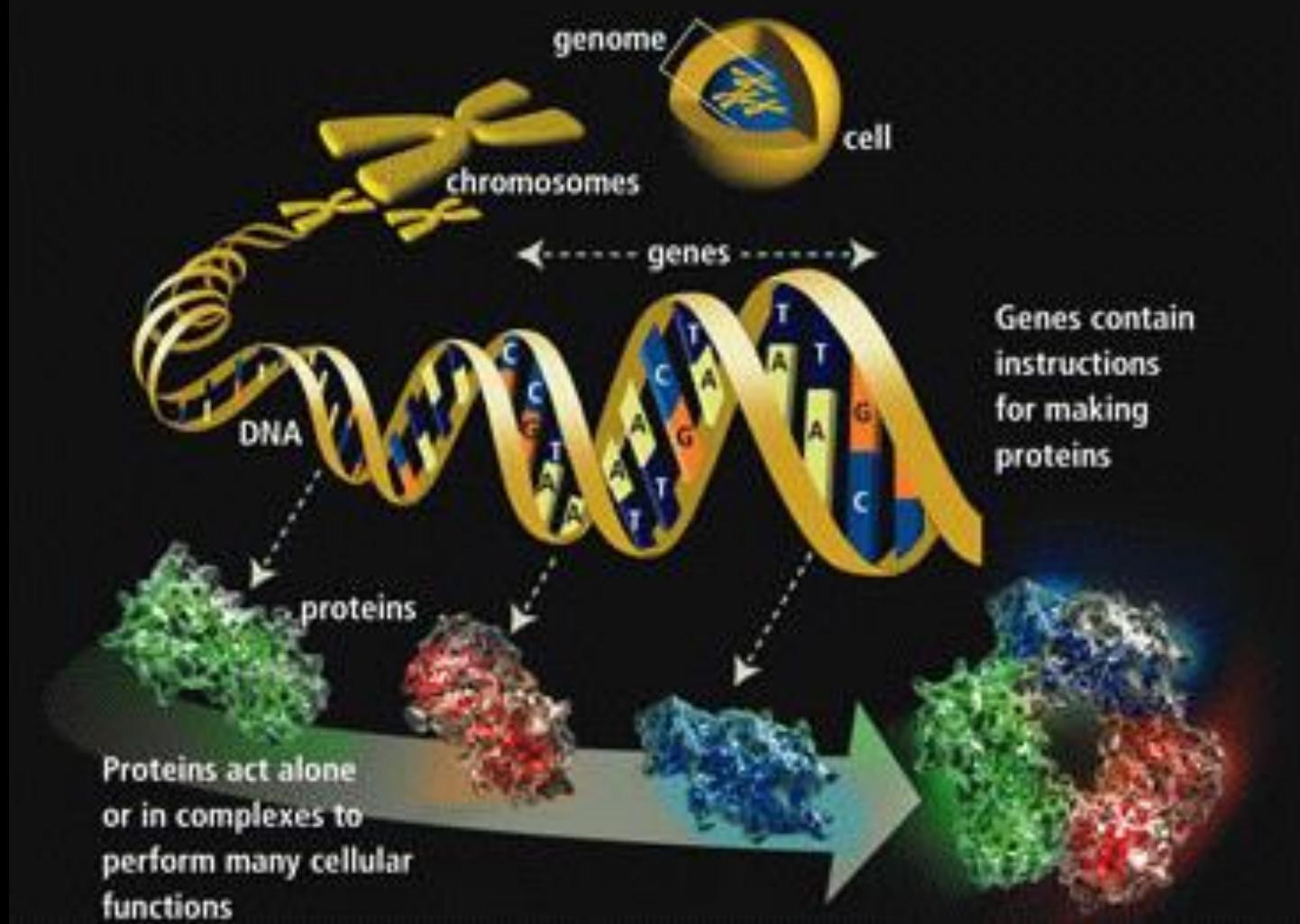


One genome -> different functions

---



# Gene expression : from genes to proteins



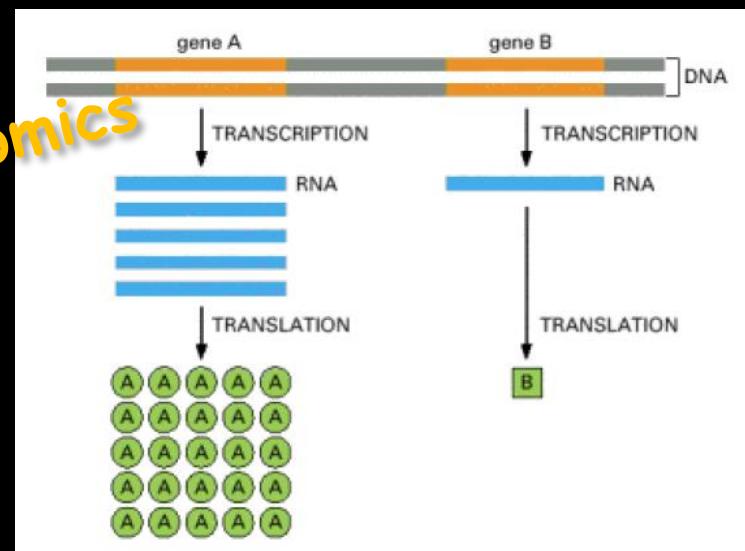
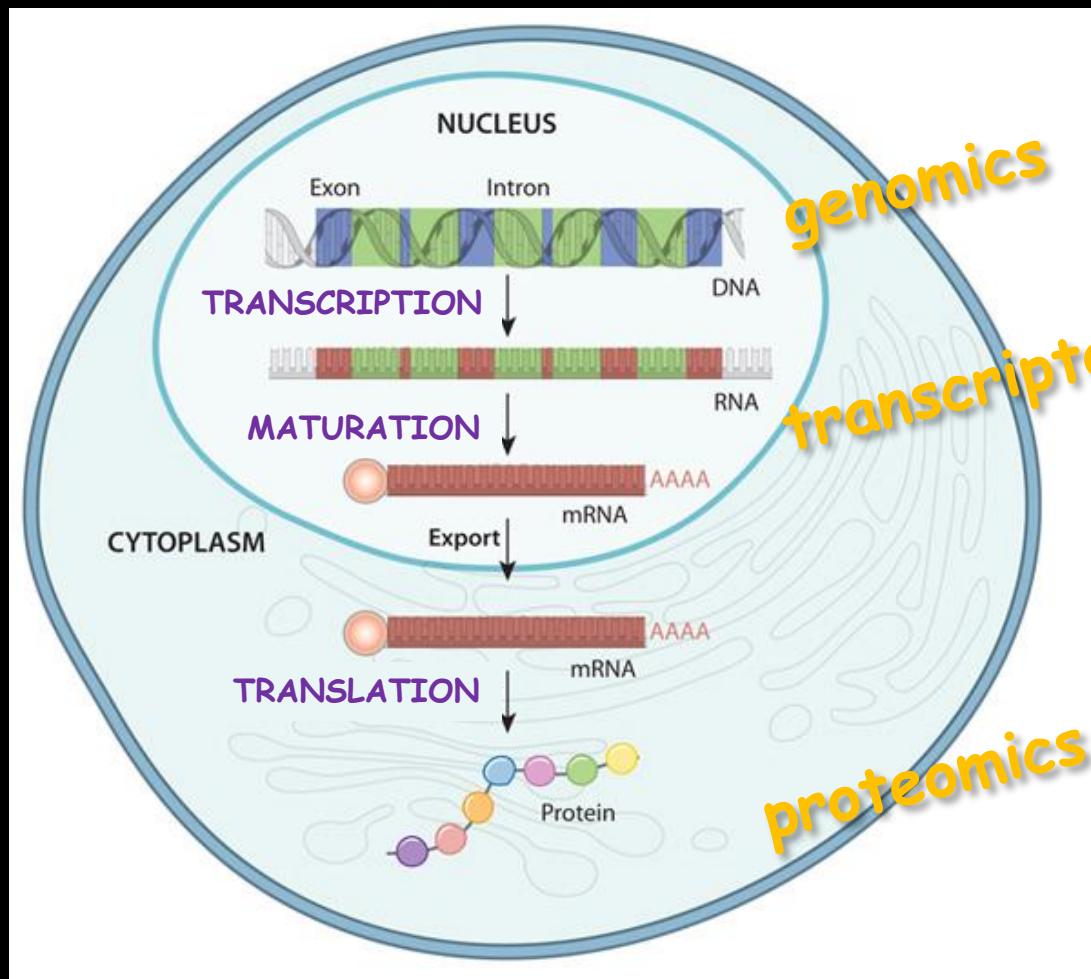
- *Gene expression*
- *Sequencing and sequences*
- *Genome annotation*
  - Structure
  - Function
  - Evolution

---

# Gene Expression

## Basis of molecular genetics

# Gene expression : from genes to proteins

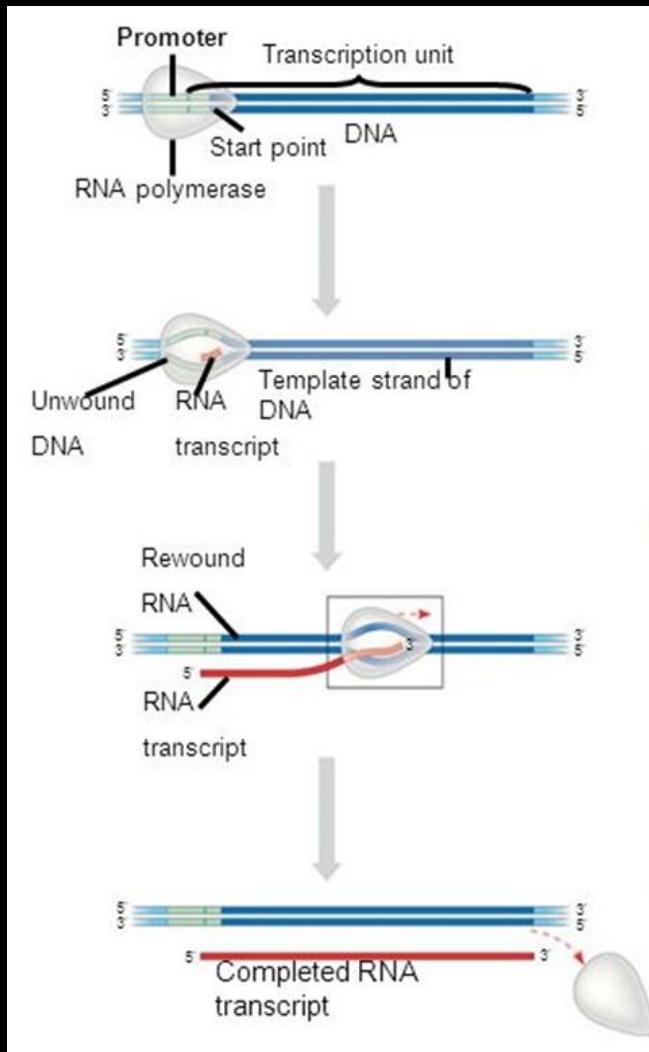


# Transcription, 3 steps

## 1. Initiation

## 2. Elongation

## 3. Termination

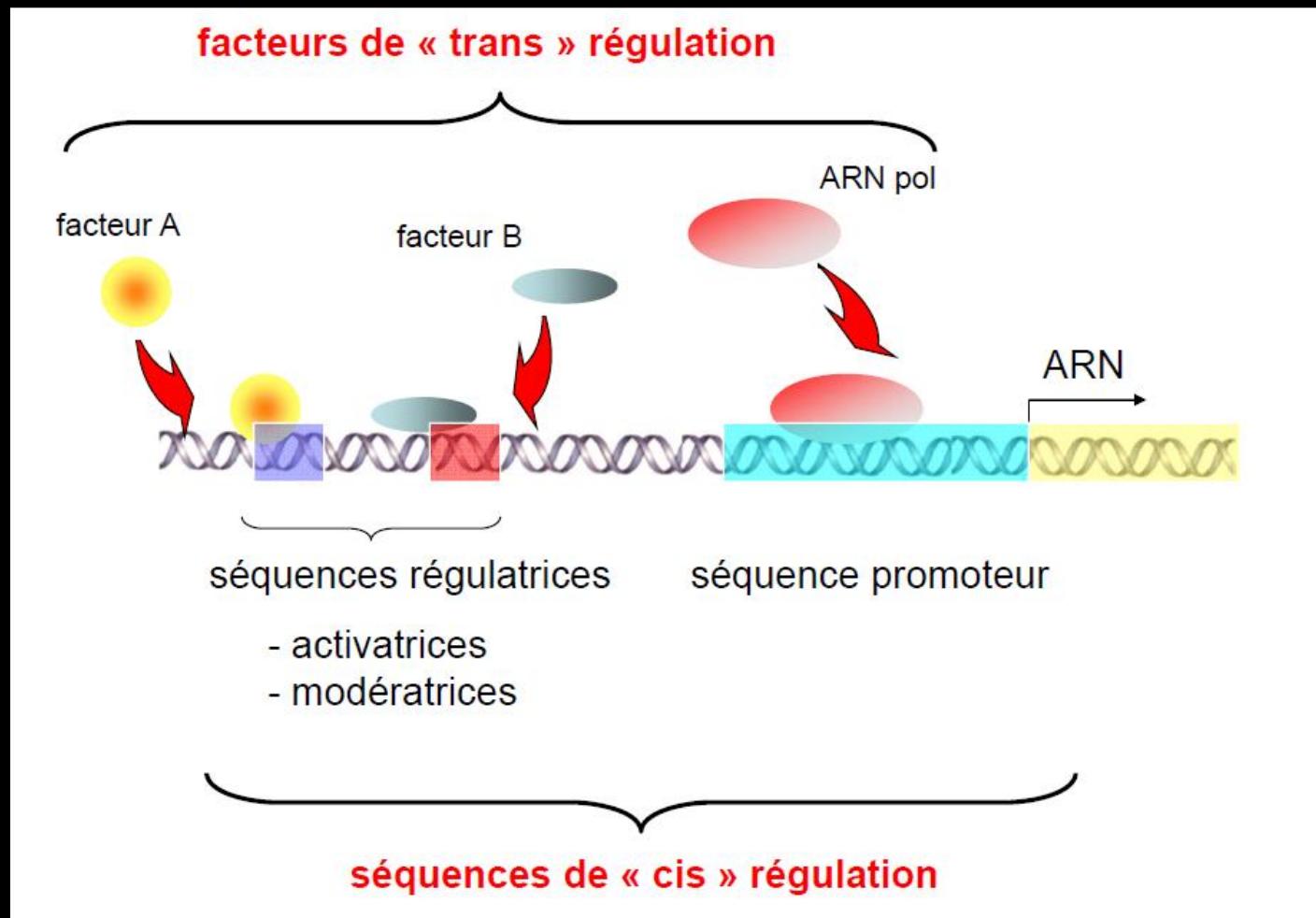


RNA polymerase binds to a promoter sequence

mRNA copy of gene (transcriptional unit) is synthesized 5' to 3'

Termination sequence causes transcription to stop

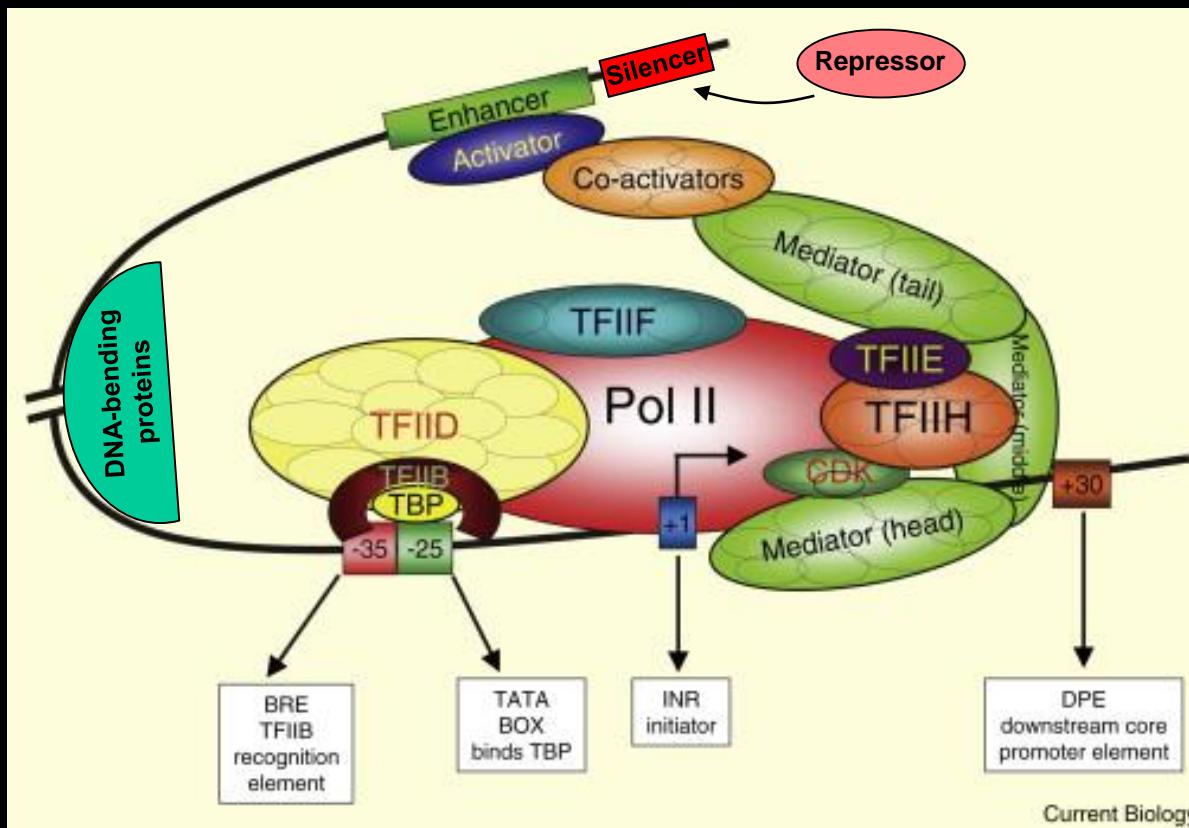
# Transcription initiation



# Transcription initiation

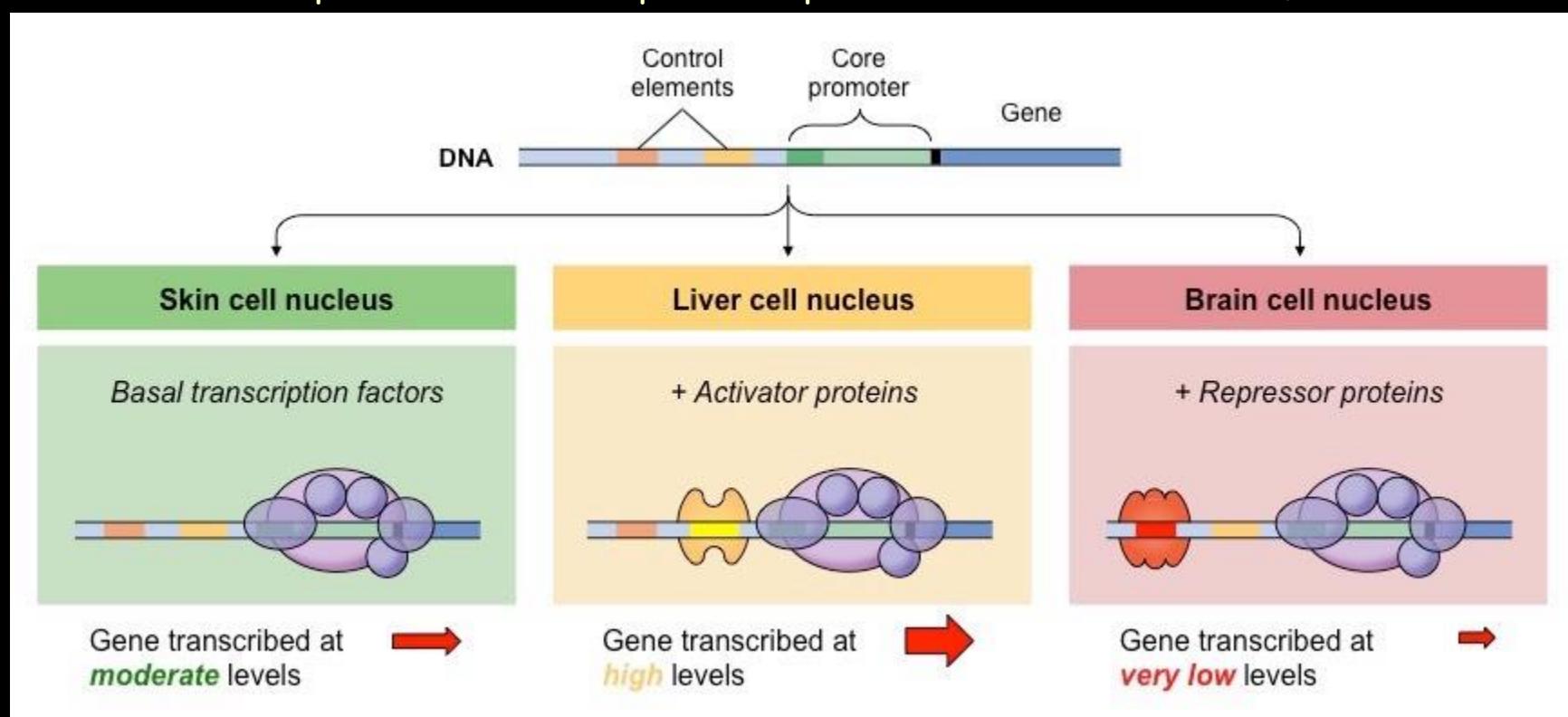
Chromatin remodeling allows DNA-protein interactions.

Transcription factors (TF), activators, co-activators, mediators recruit RNA polymerase II to constitute the transcription preinitiation complex.

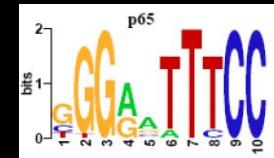


~ 100 proteins

# Transcription regulation

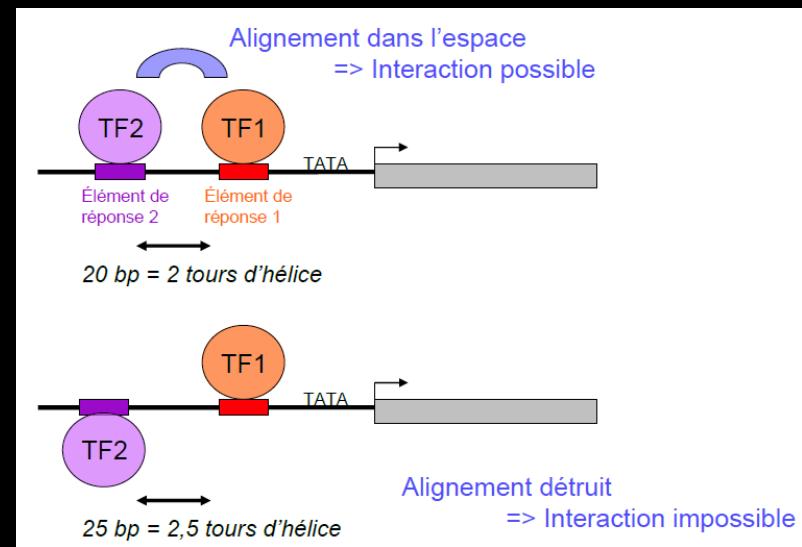
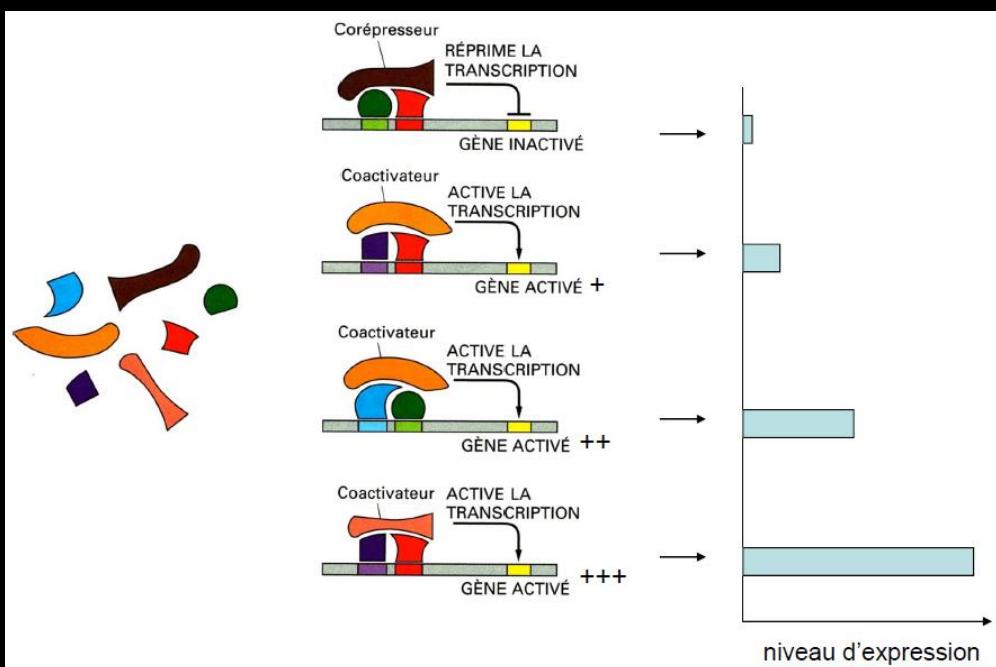
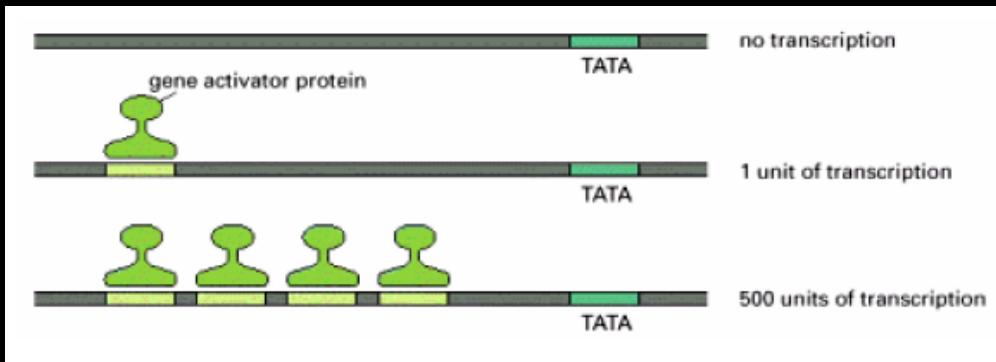


Each TF recognizes a target sequence, TFBS, such as :



# Transcription regulation

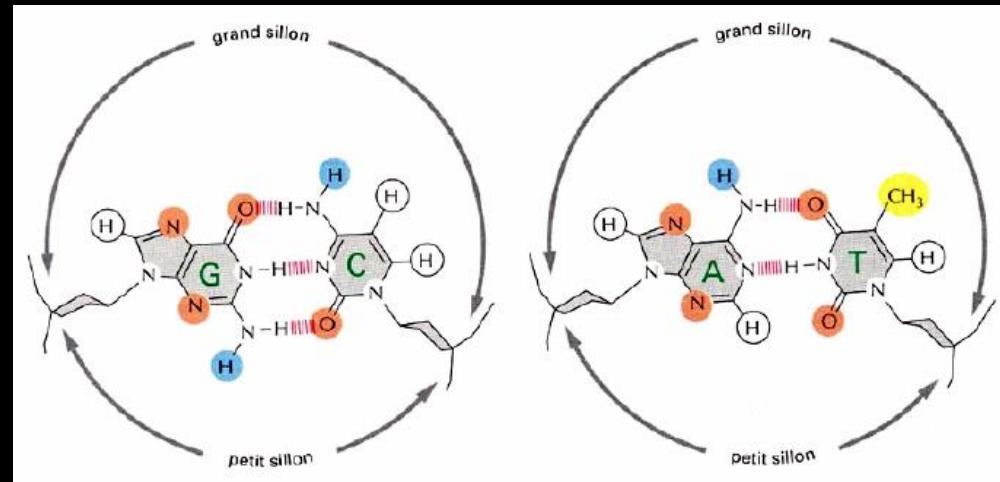
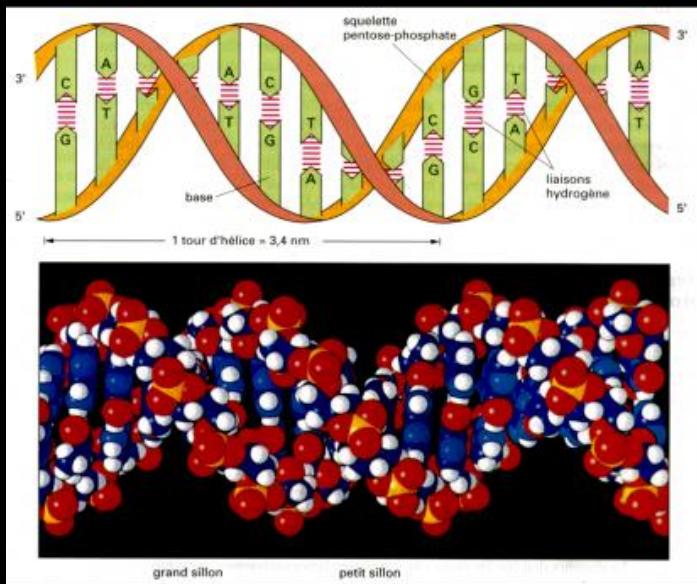
## Cooperation/competition between TFs for regulation



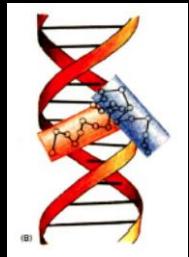
► Applications in biotechnologies for modulate/control gene expression

# DNA-protein interactions

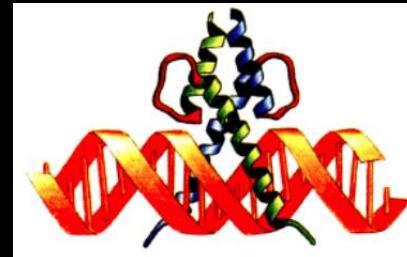
DNA-binding proteins mainly interact with the major groove of DNA



TFs are classified according their secondary structure which defined the interaction with double helix : basic-helix-loop-helix (leucine zipper...), zinc-fingers, helix-turn-helix, beta scaffold with minor groove contact...

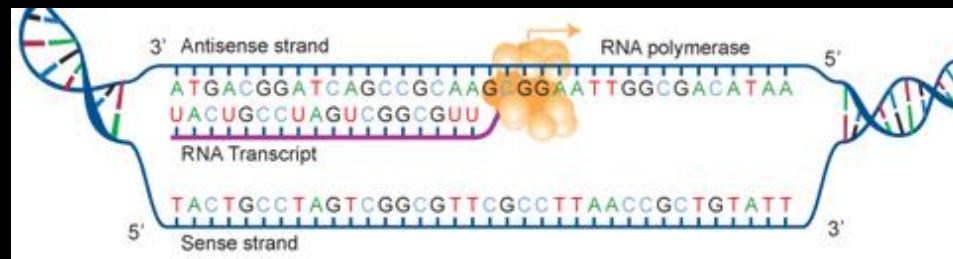
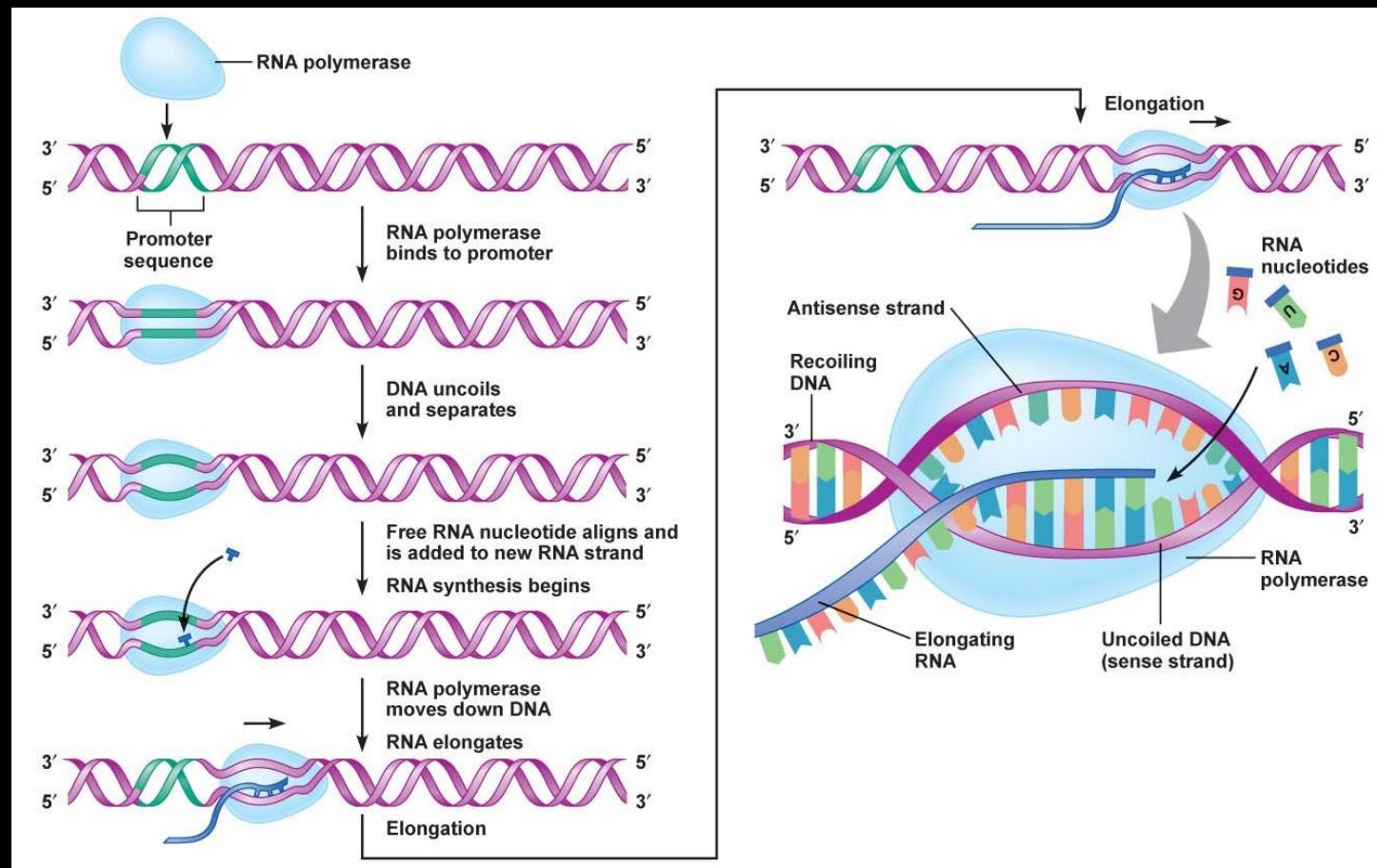


helix-turn-helix

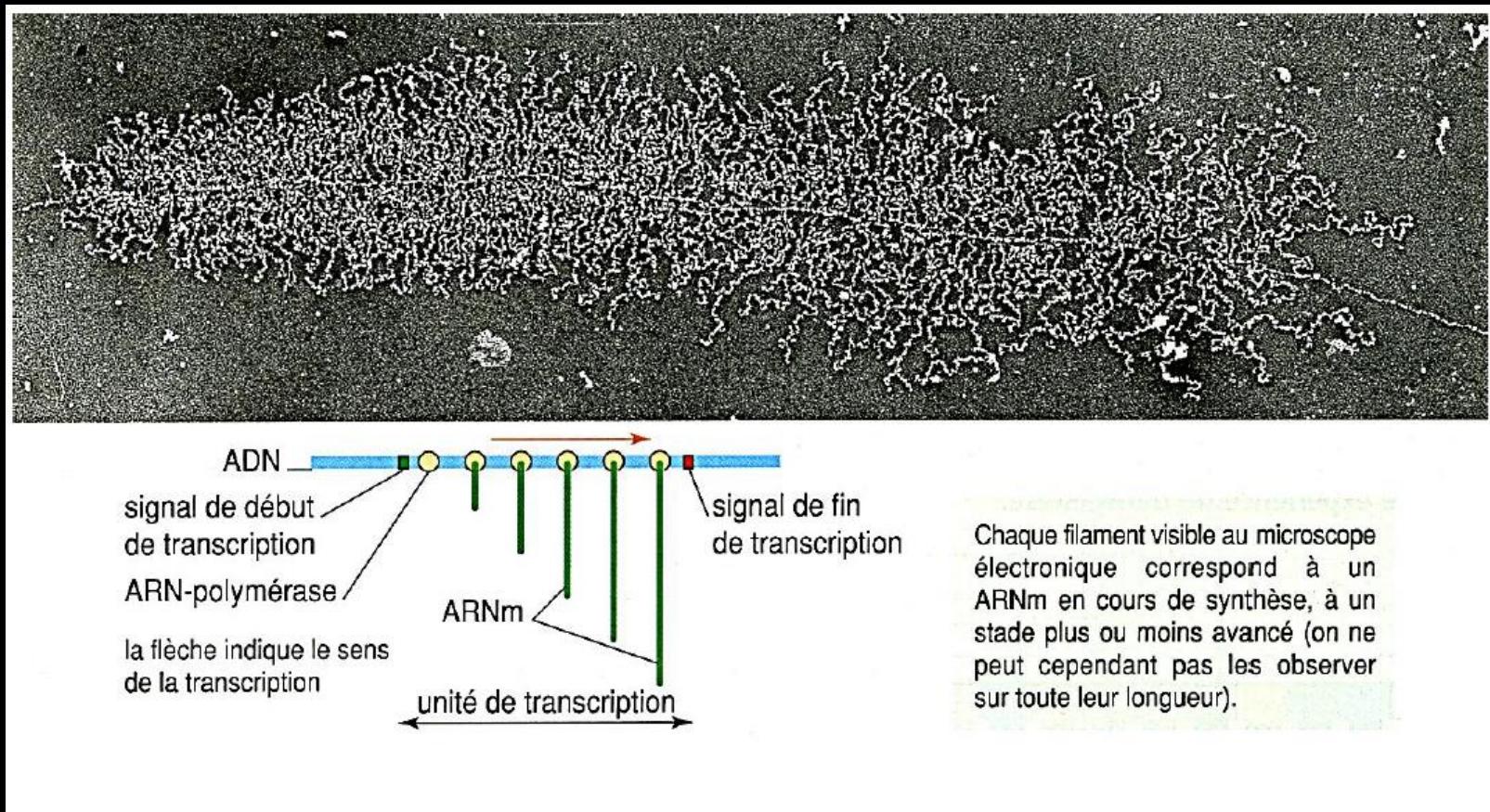


dimeric  
helix-loop-helix

# Transcription elongation

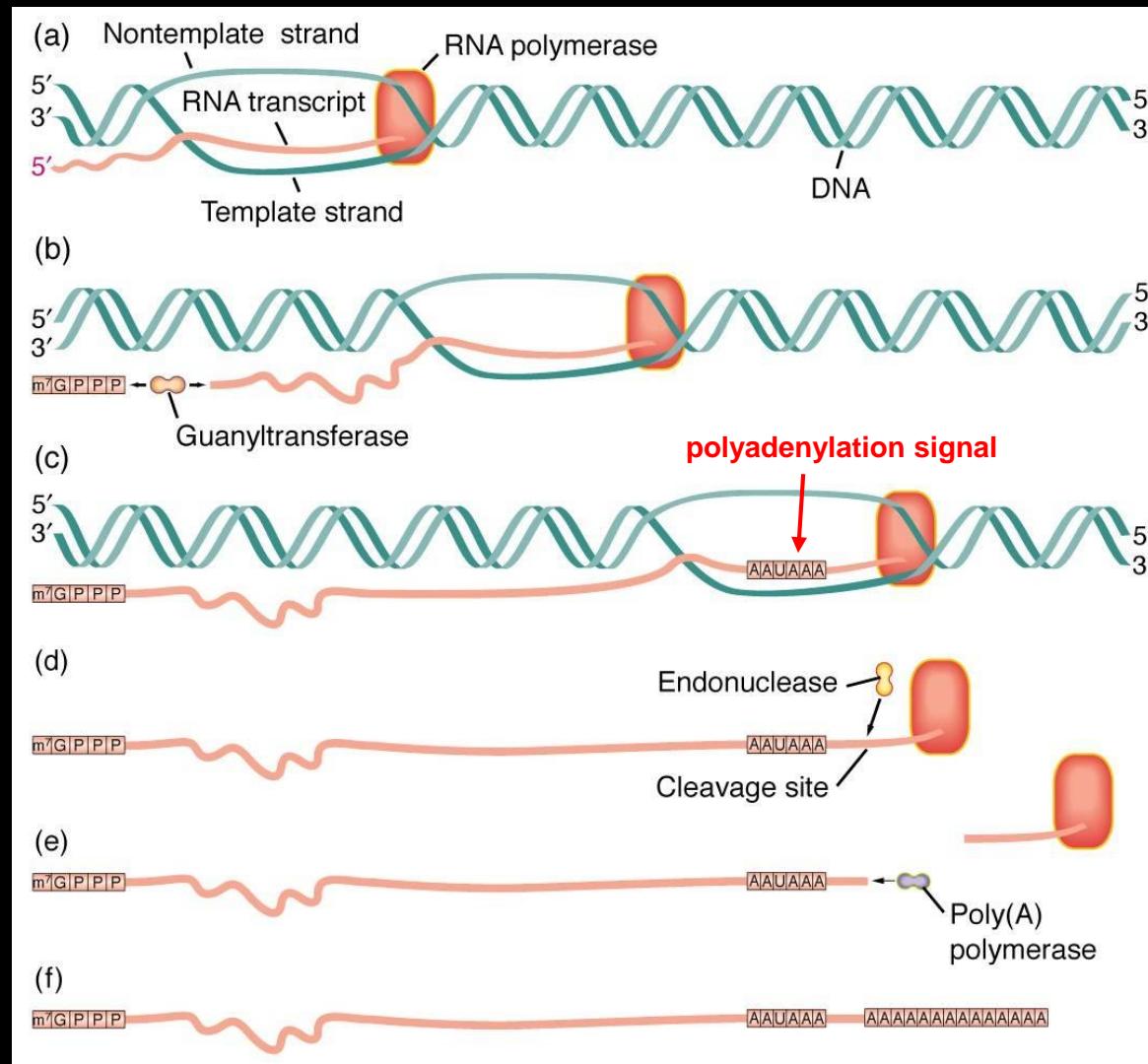


# Transcription elongation



~ 40-80 nucleotides / sec

# Transcription termination and maturation

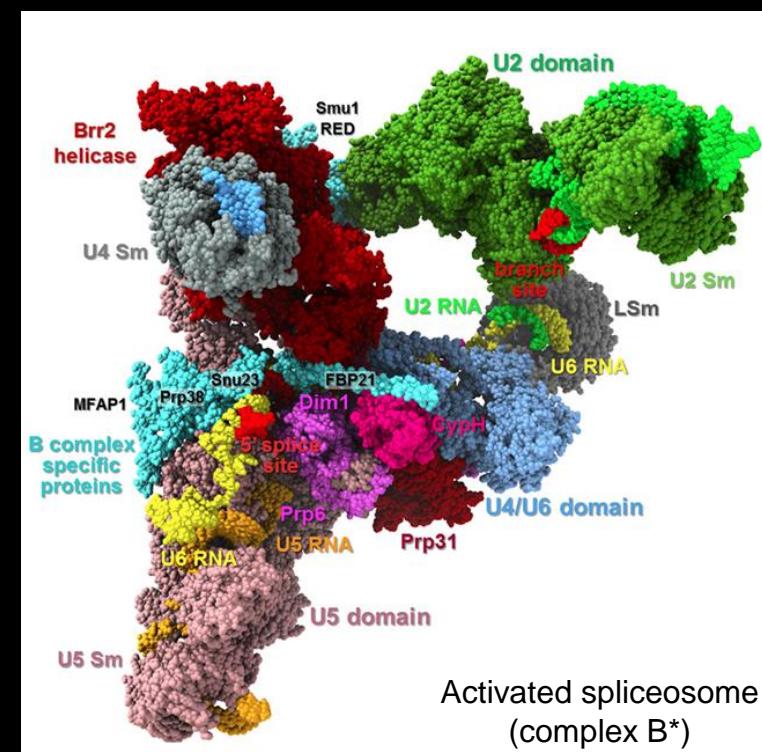
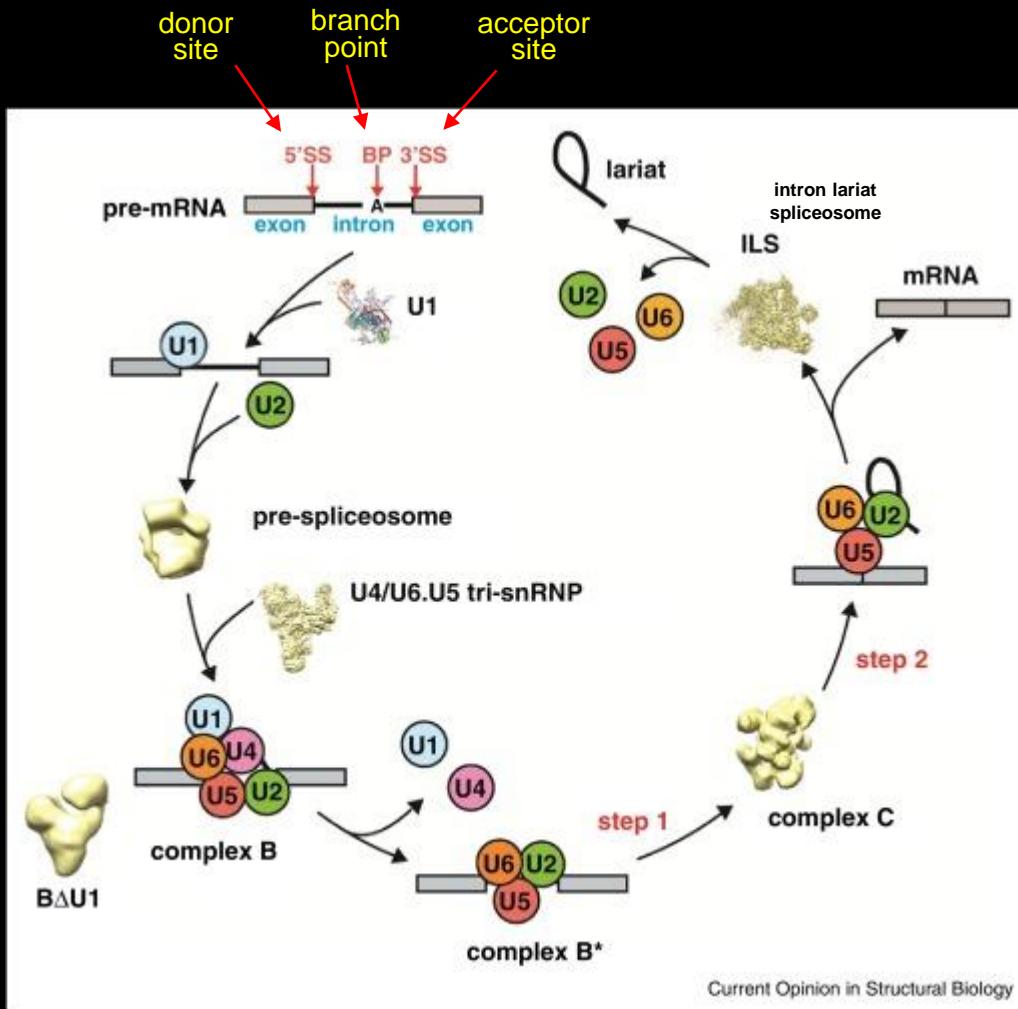


cap in 5'

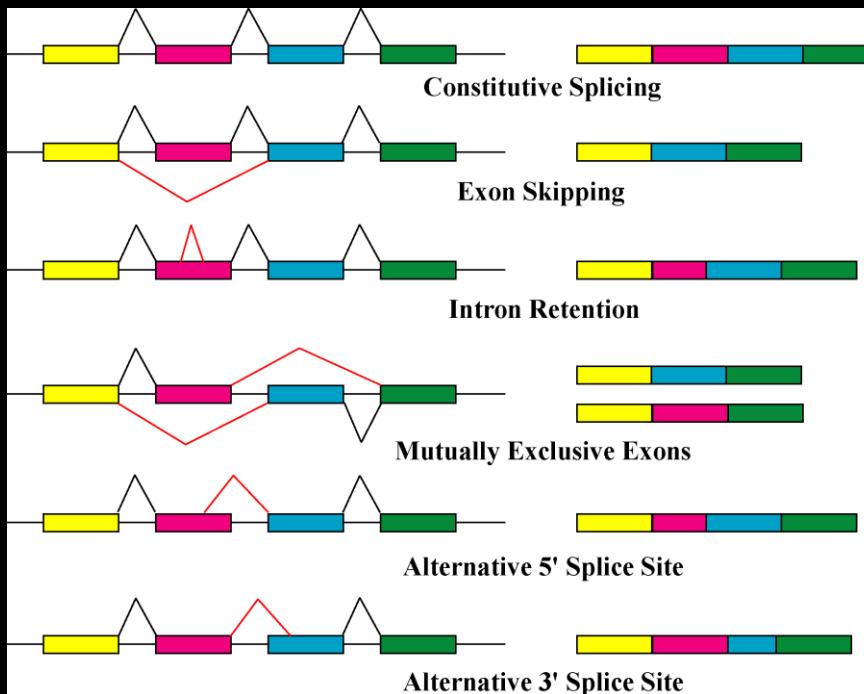
Poly(A) tail in 3'

# Pre-mRNA splicing : intron removing

Small nuclear ribonucleic acids (snRNA, ~150 b) and proteins constitute small nuclear ribonucleoproteins (snRNP), components of **spliceosome**.



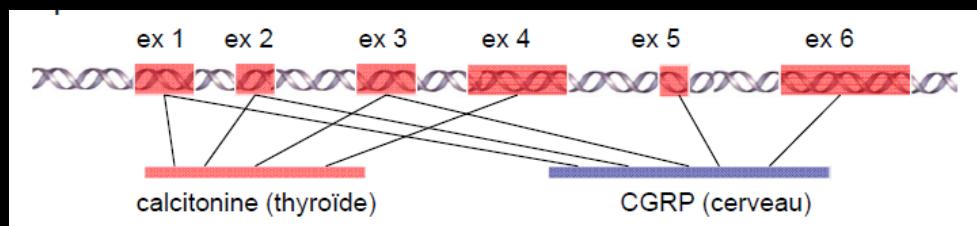
# Post-transcriptional regulation: alternative splicing



One pre-mRNA can produce a large diversity of mature RNA

Exemple 1 :

2 peptides  
2 functions



Hormone  
reduces  $\text{Ca}^{2+}$  in blood

Vasodilator  
nociception

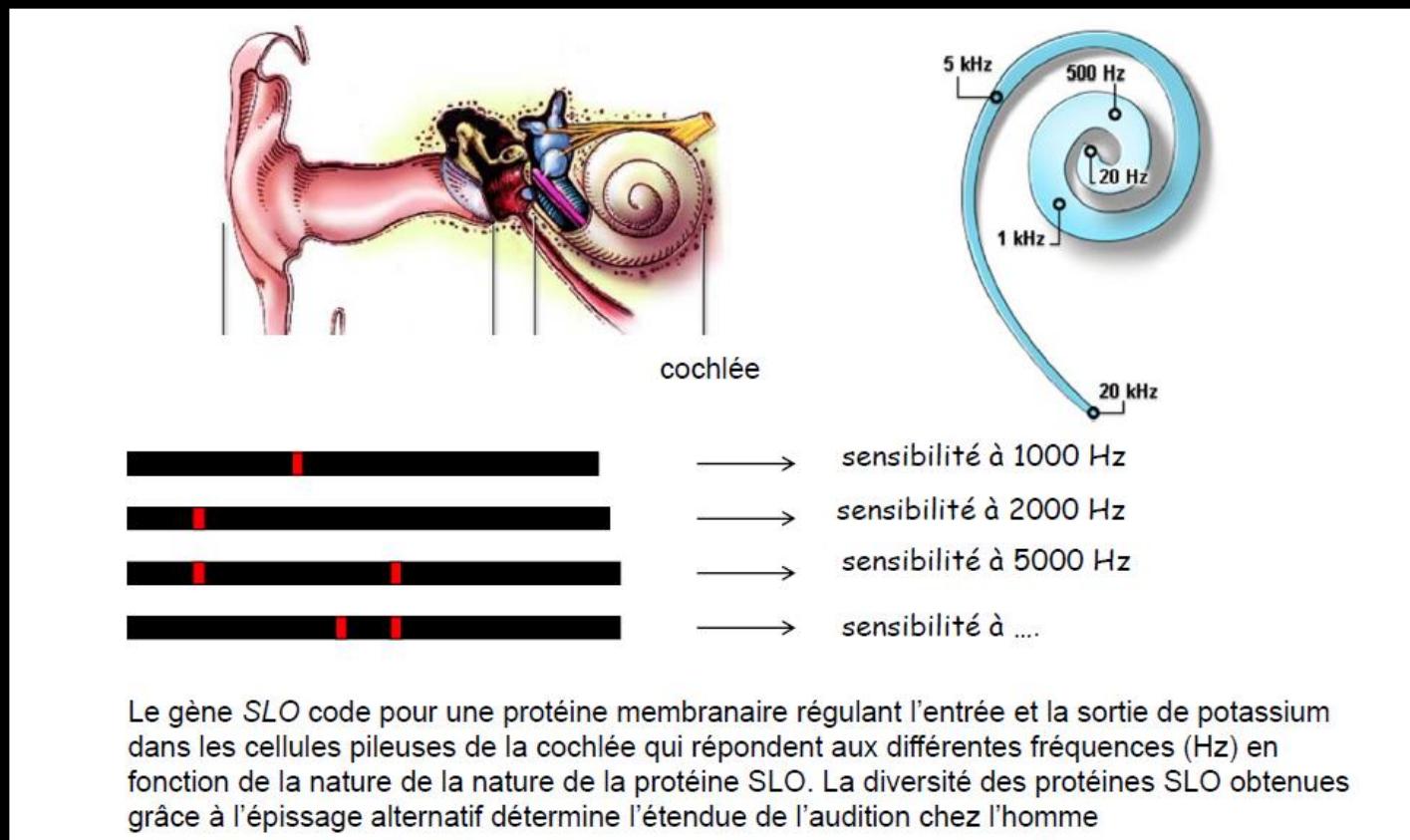
Calcitonin  
Gene-  
Related  
Peptide

# Post-transcriptional regulation: alternative splicing

Exemple 2 :



Gene *slo* : 35 exons including 8 optional ones  
500 different transcripts ► 500 proteins SLO

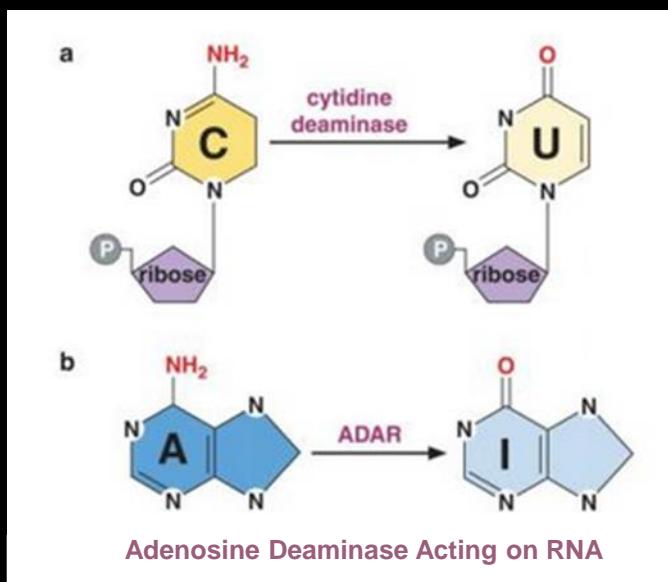


# Post-transcriptional regulation: RNA editing

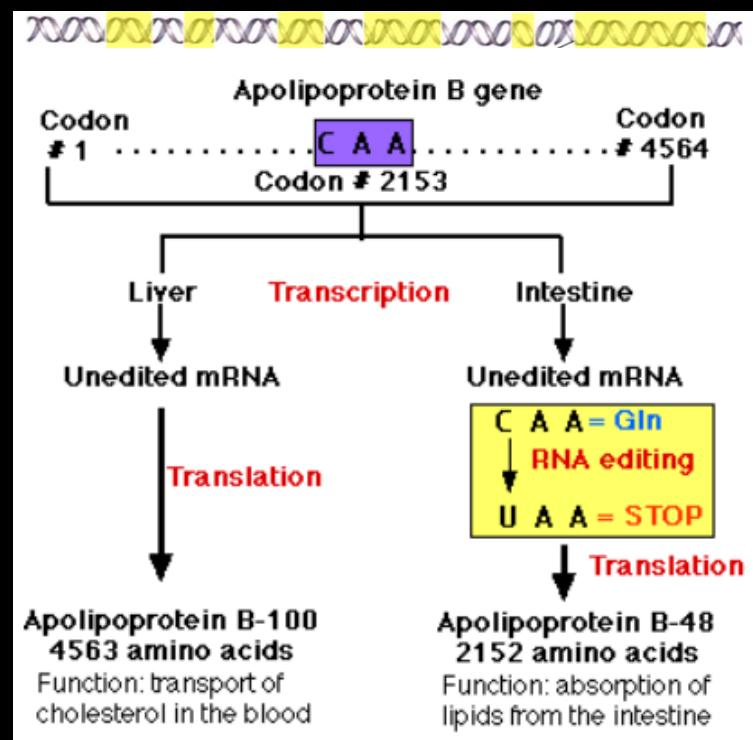
RNA editing is a dynamic mechanism that generates site-specific deamination. The most common type of RNA editing involves nucleotide substitution that consists of :

Adenosine A  $\longrightarrow$  inosine I (mainly in human), I is read as G by the translation machinery of the cell.

Cytidine C  $\longrightarrow$  uridine U (mainly in plants).

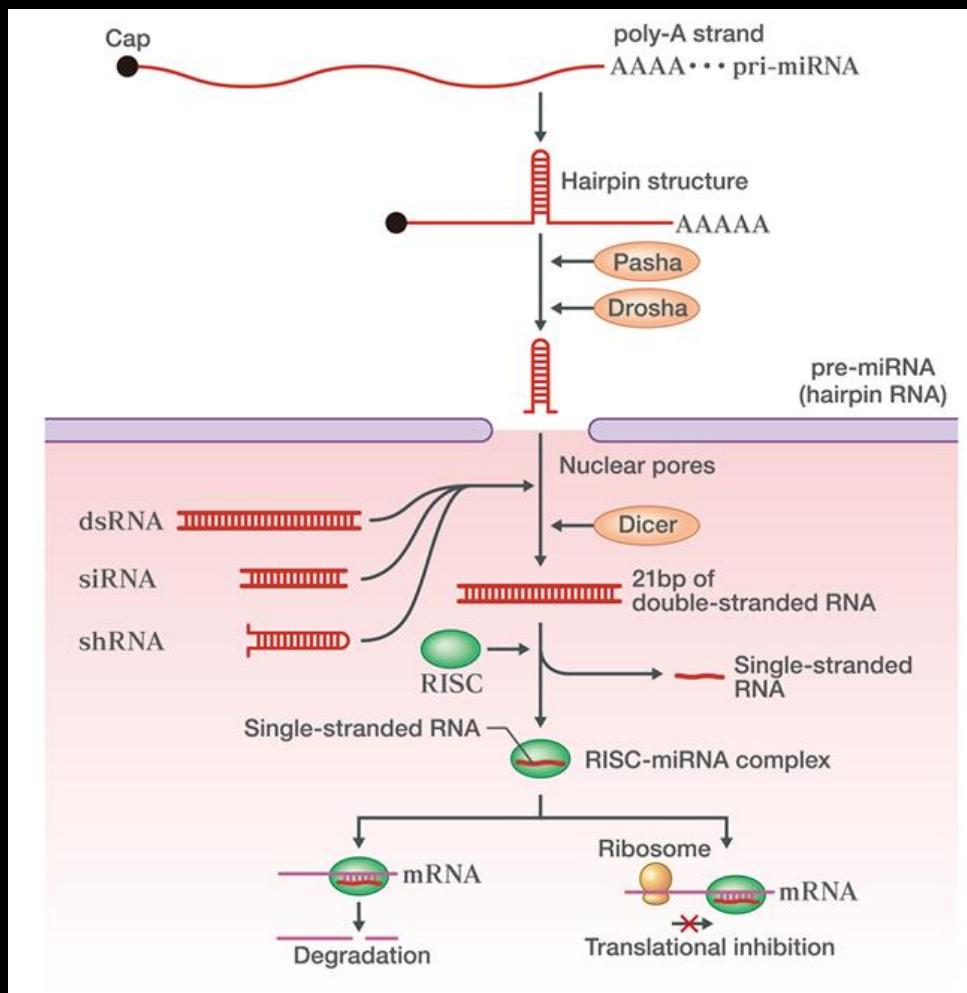


Exemple

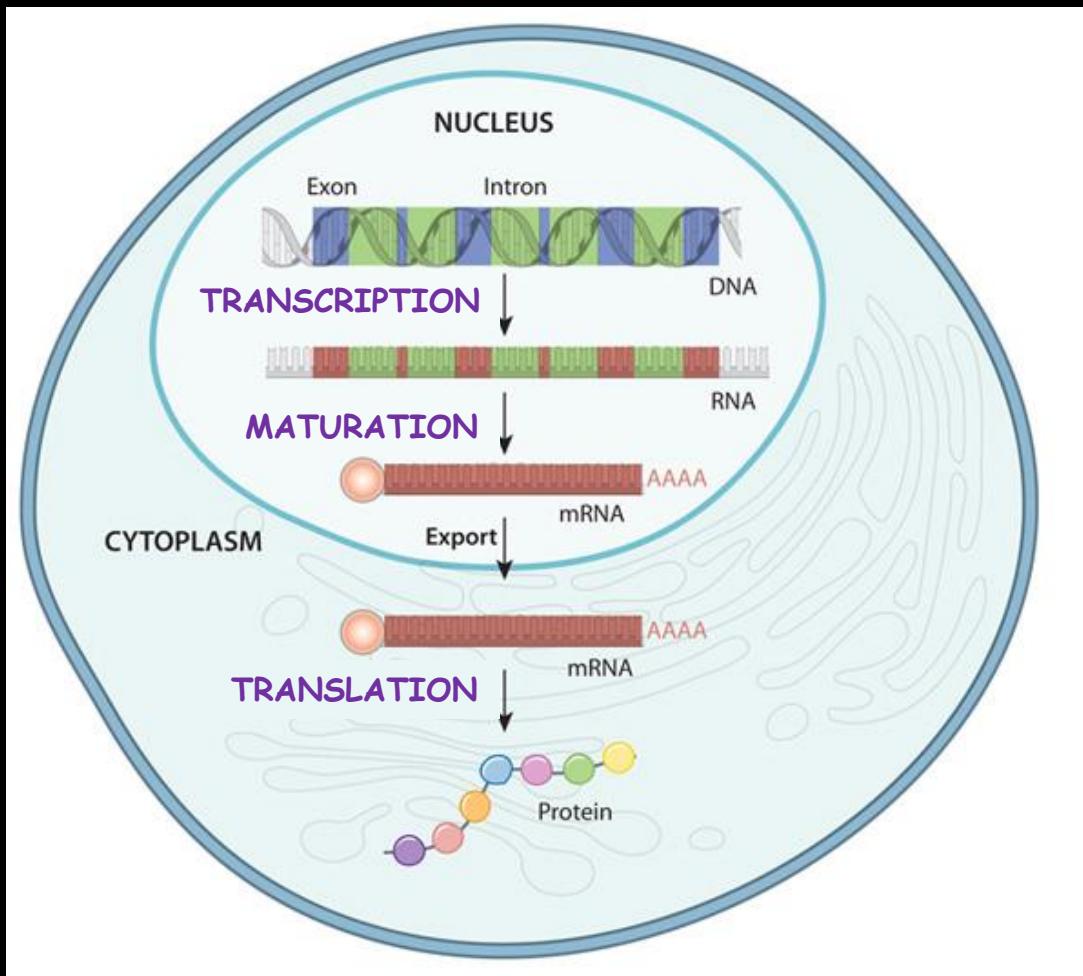


# Post-transcriptional regulation: miRNA

MicroRNAs (miRNAs) are a class of small non-coding RNAs that post-transcriptionally control gene expression via either translational repression or mRNA degradation.



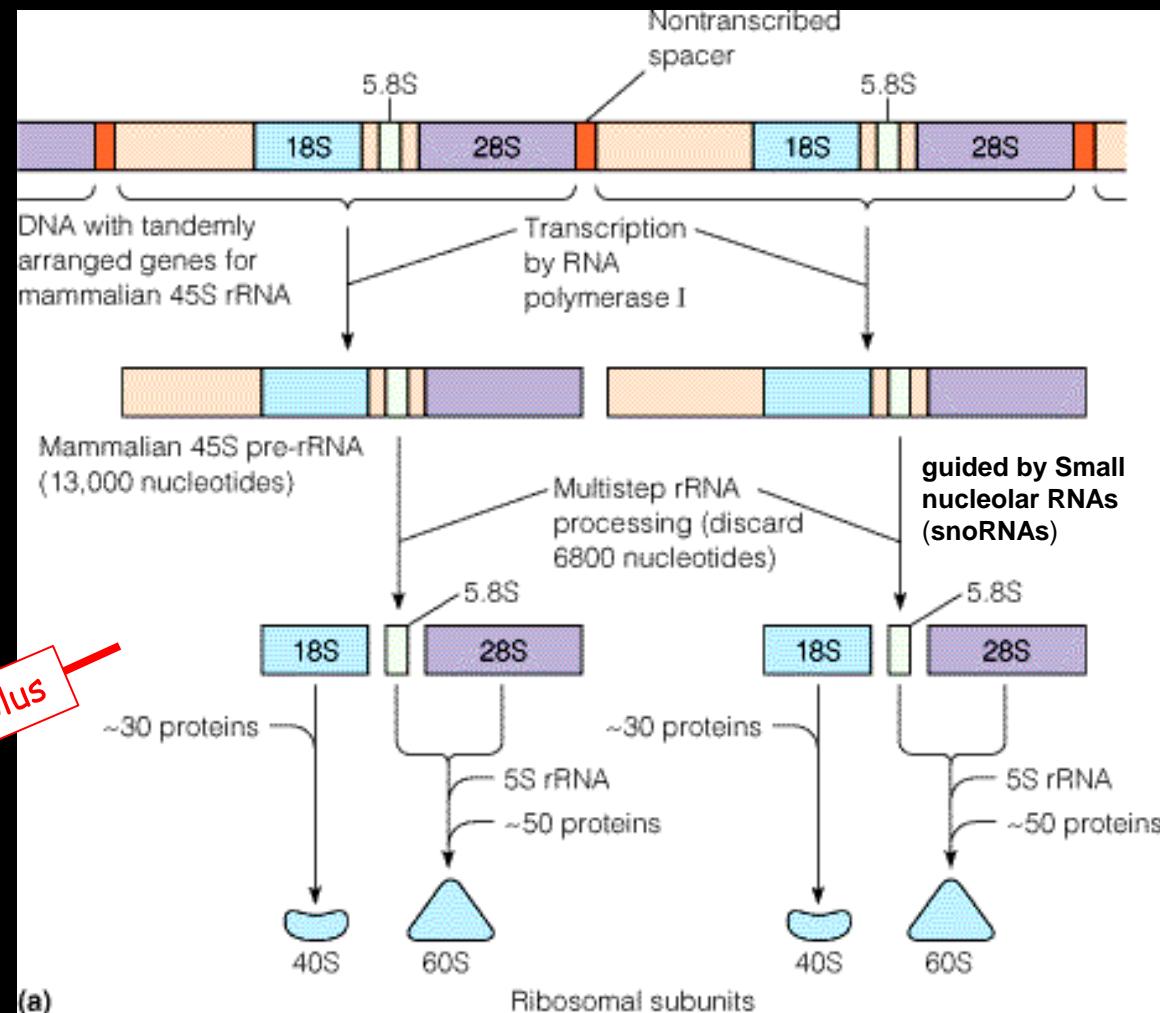
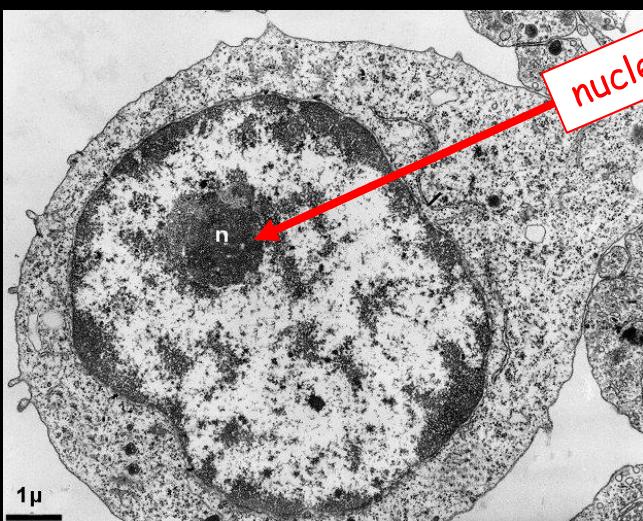
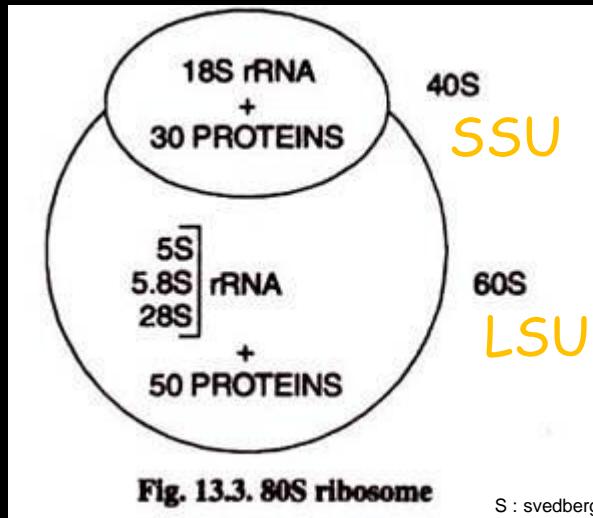
# Translation



Actors  
mature mRNA  
Ribosomes  
tRNAs

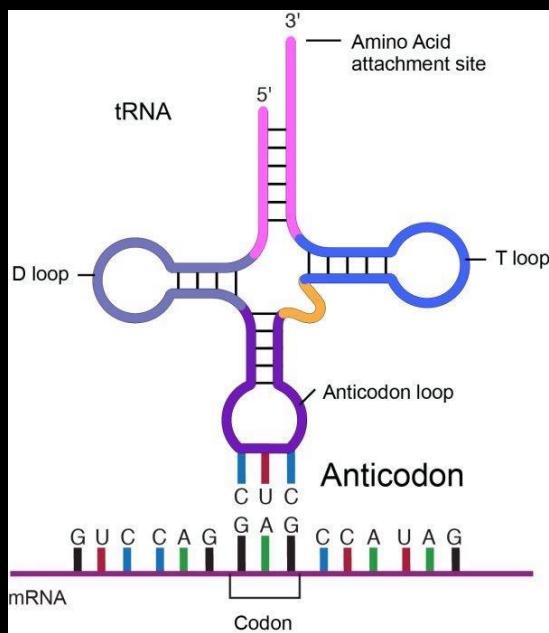
Place  
Cytoplasm  
/  
ER associated

# Ribosome : translational machinery

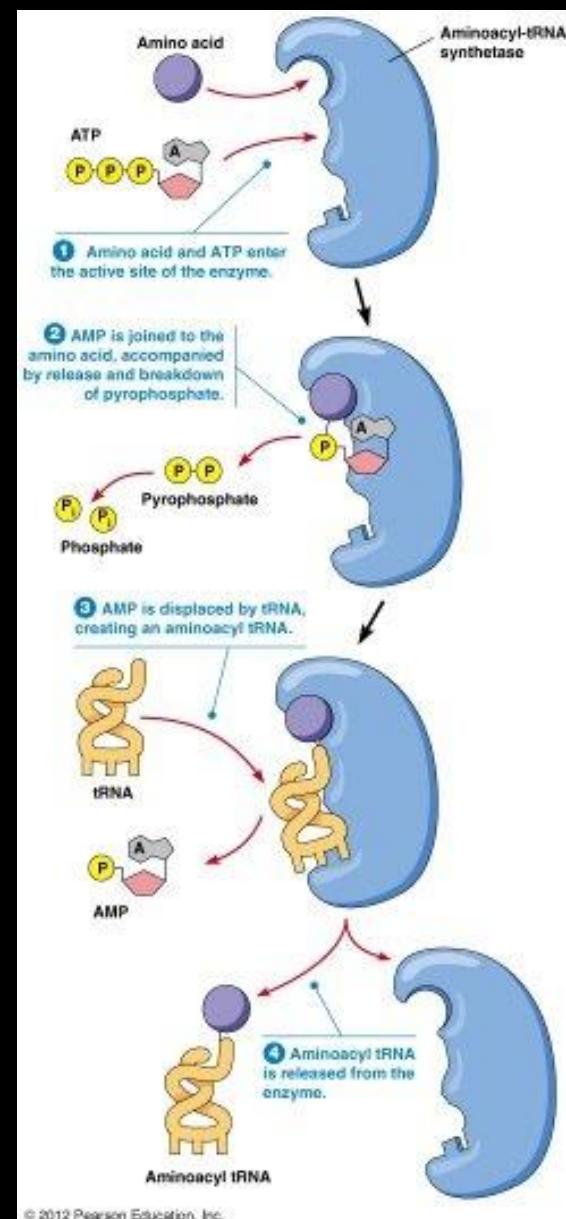


rRNA : ~90% of the RNA in a cell

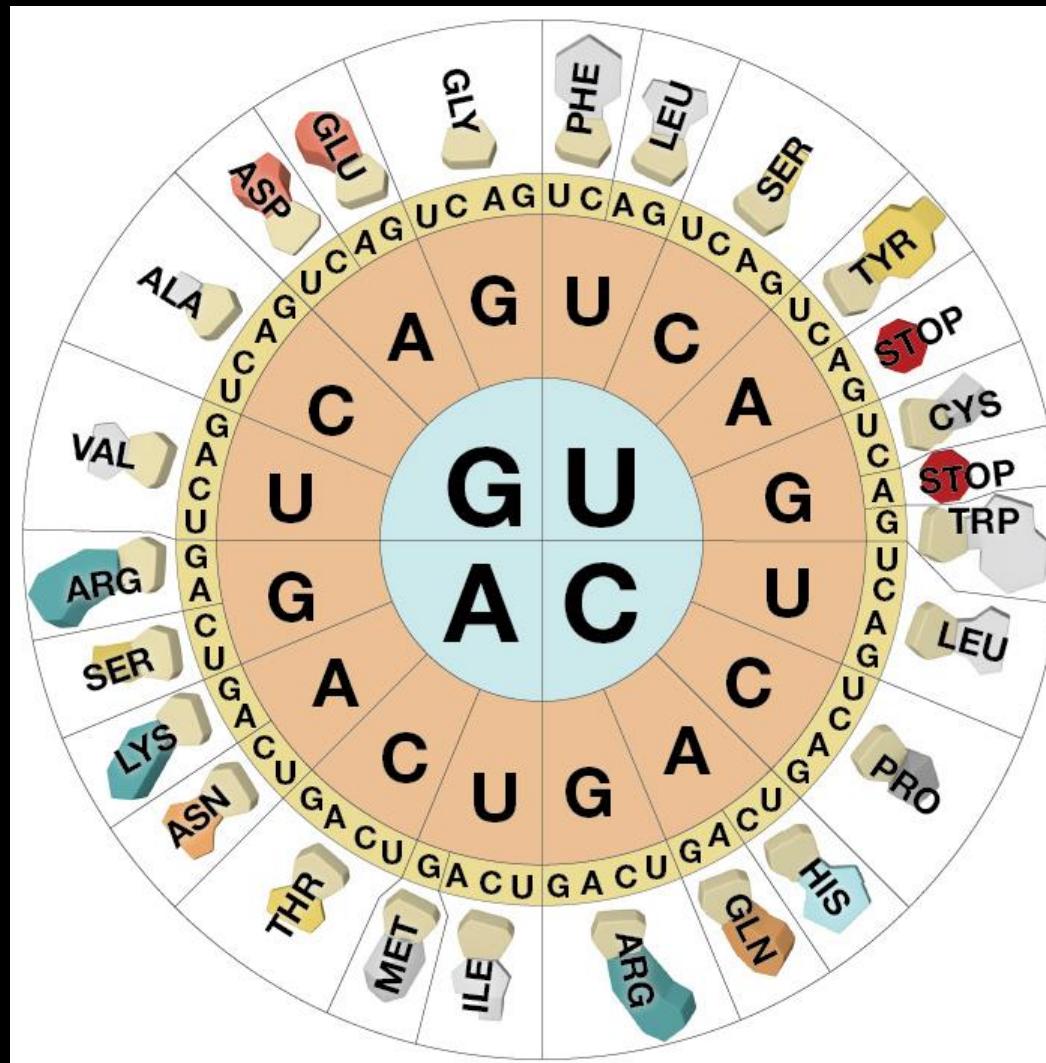
# tRNA : links between nucleotides and amino-acids



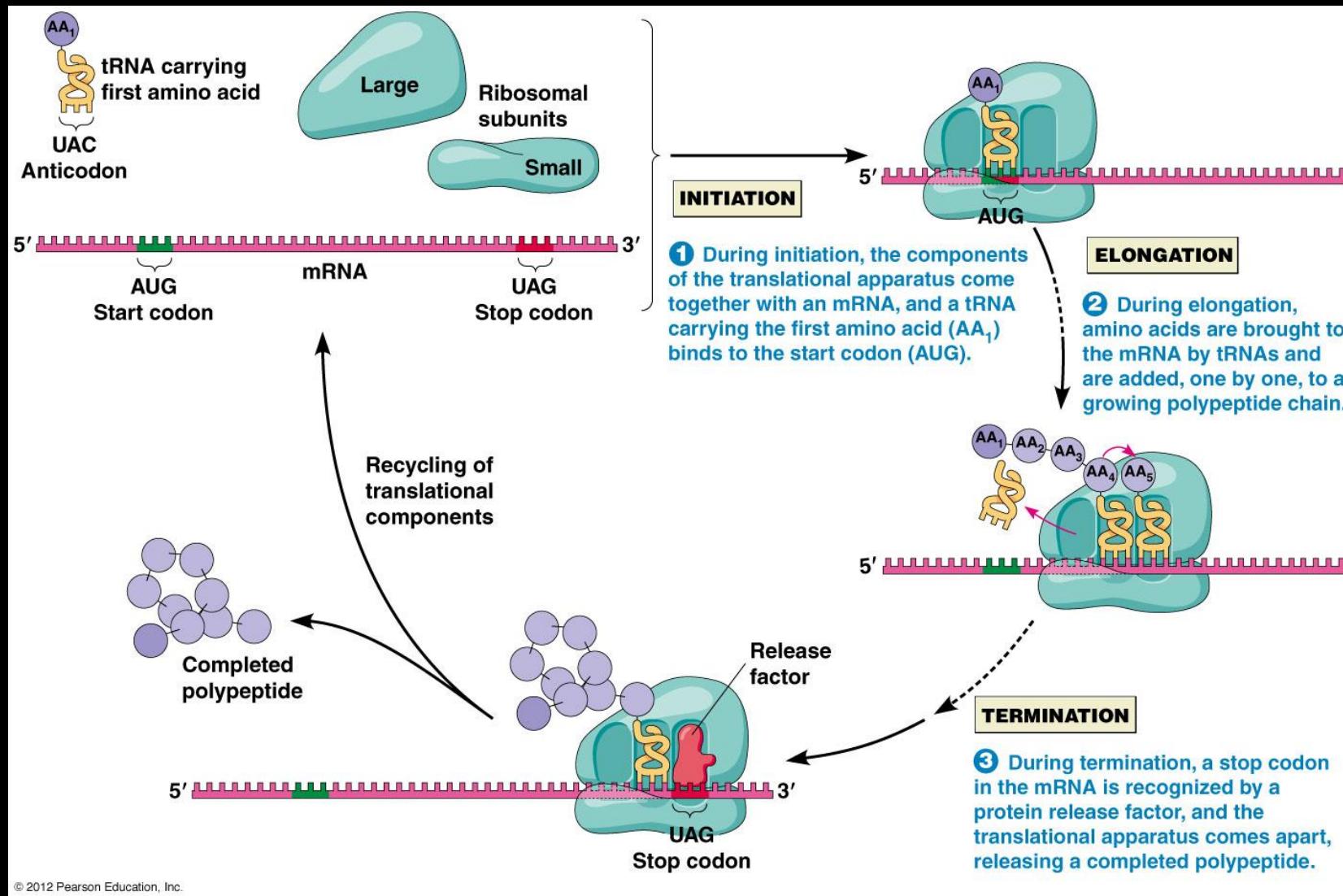
Stem loop structure  
One tRNA per codon



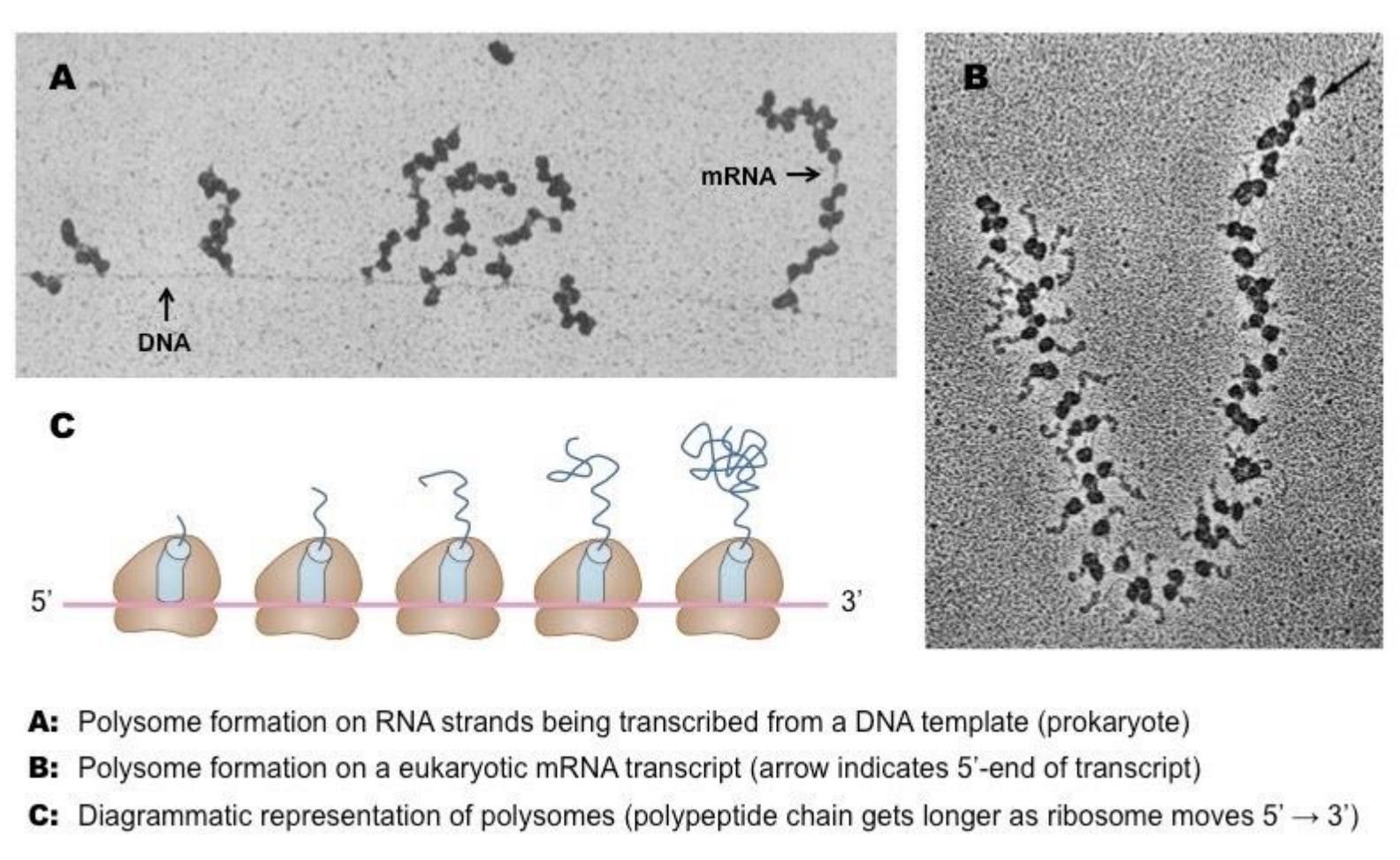
# The genetic code



# Translation

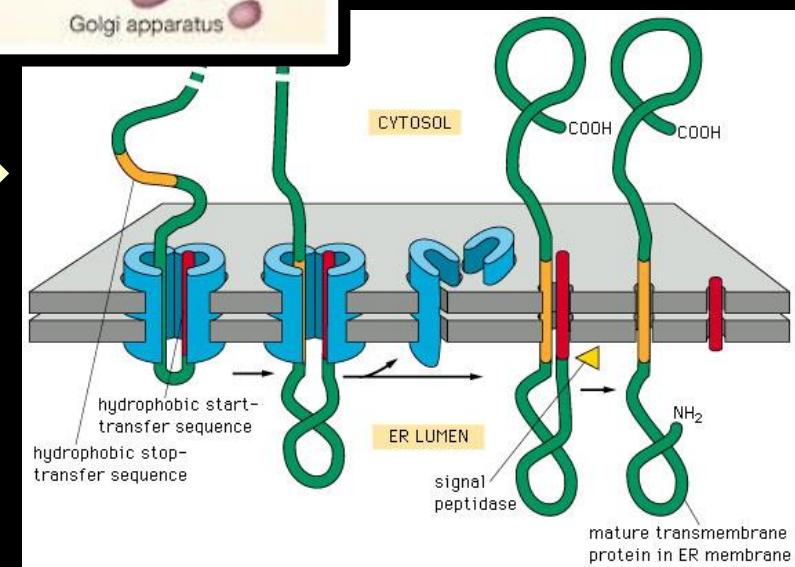
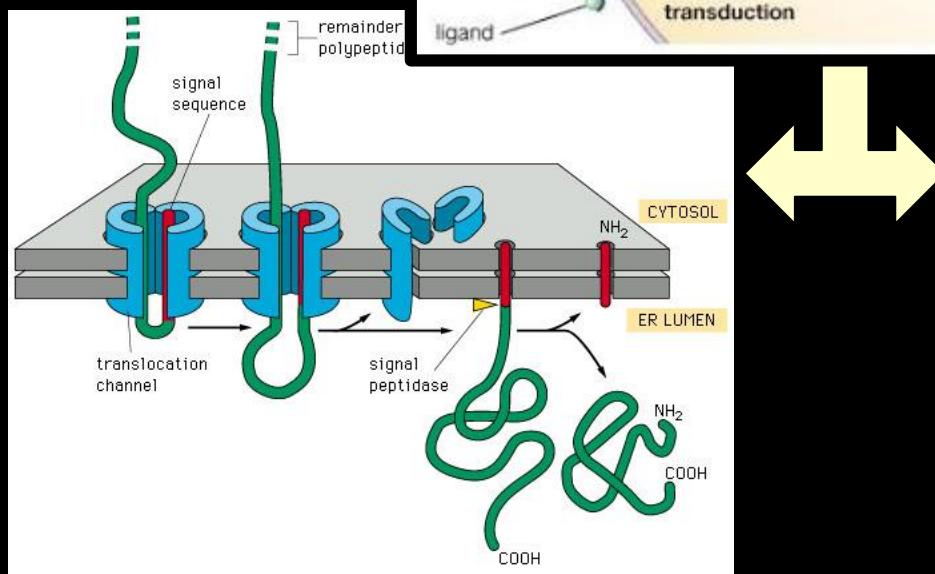
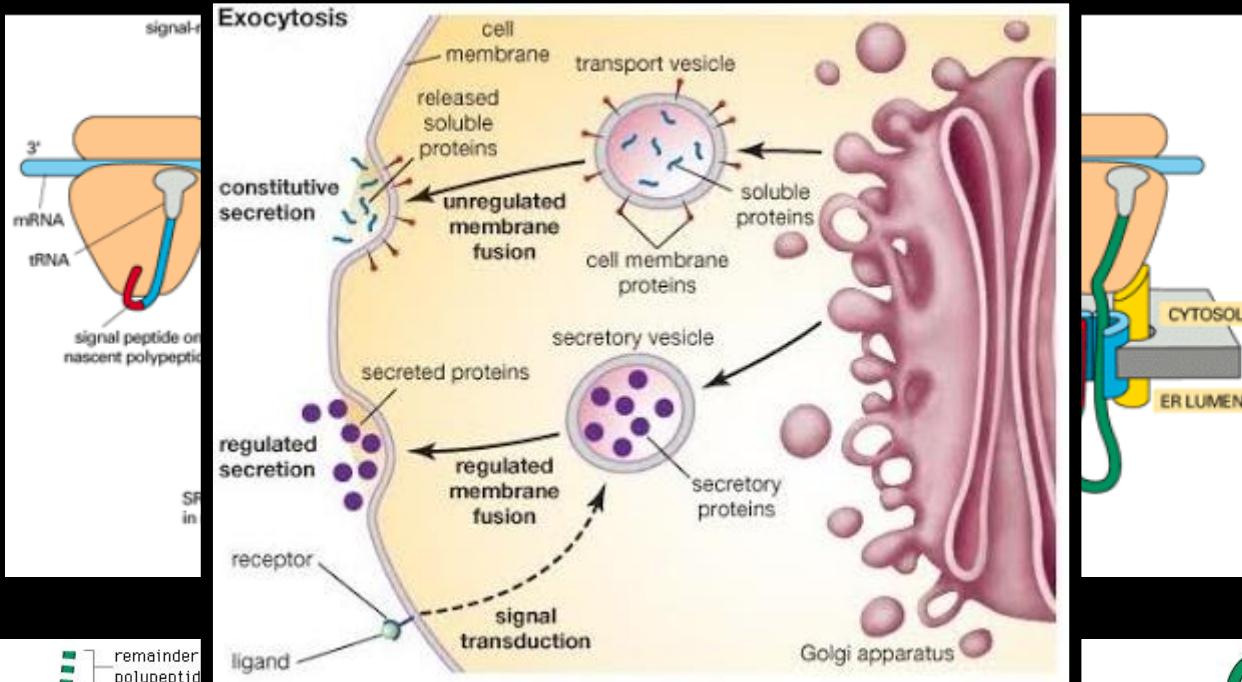


# Polysomes

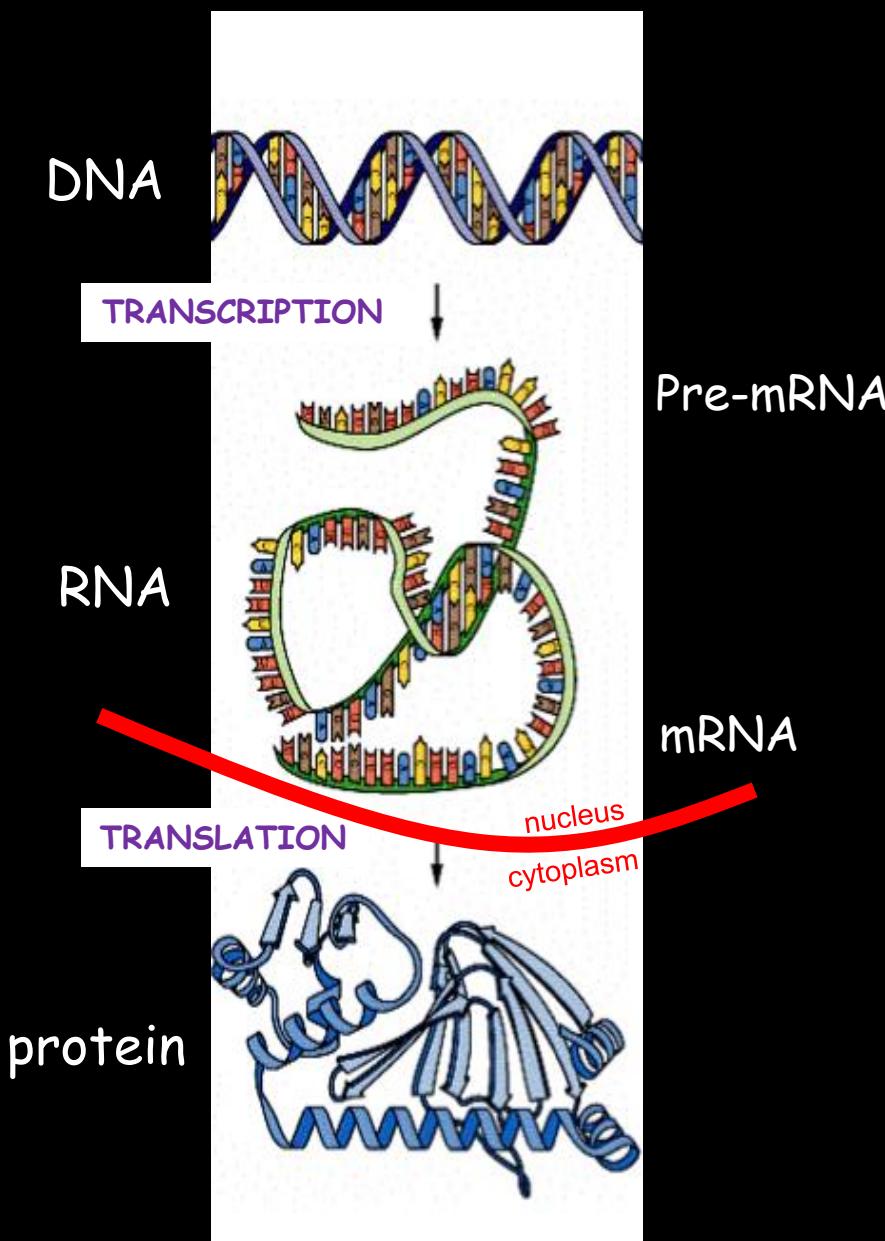


~ 5 aa (eukaryote) to 50 aa (prokaryote) / sec

# Proteins targeted to the endoplasmic reticulum



# Levels of regulation of gene expression



**Transcriptional control**  
DNA access  
chromatin structure  
DNA methylation  
Transcription factors binding  
Initiation, elongation, termination

**Post-transcriptional control**  
Maturation  
Splicing  
Editing  
Stabilization (miRNA control, polyA)  
Export

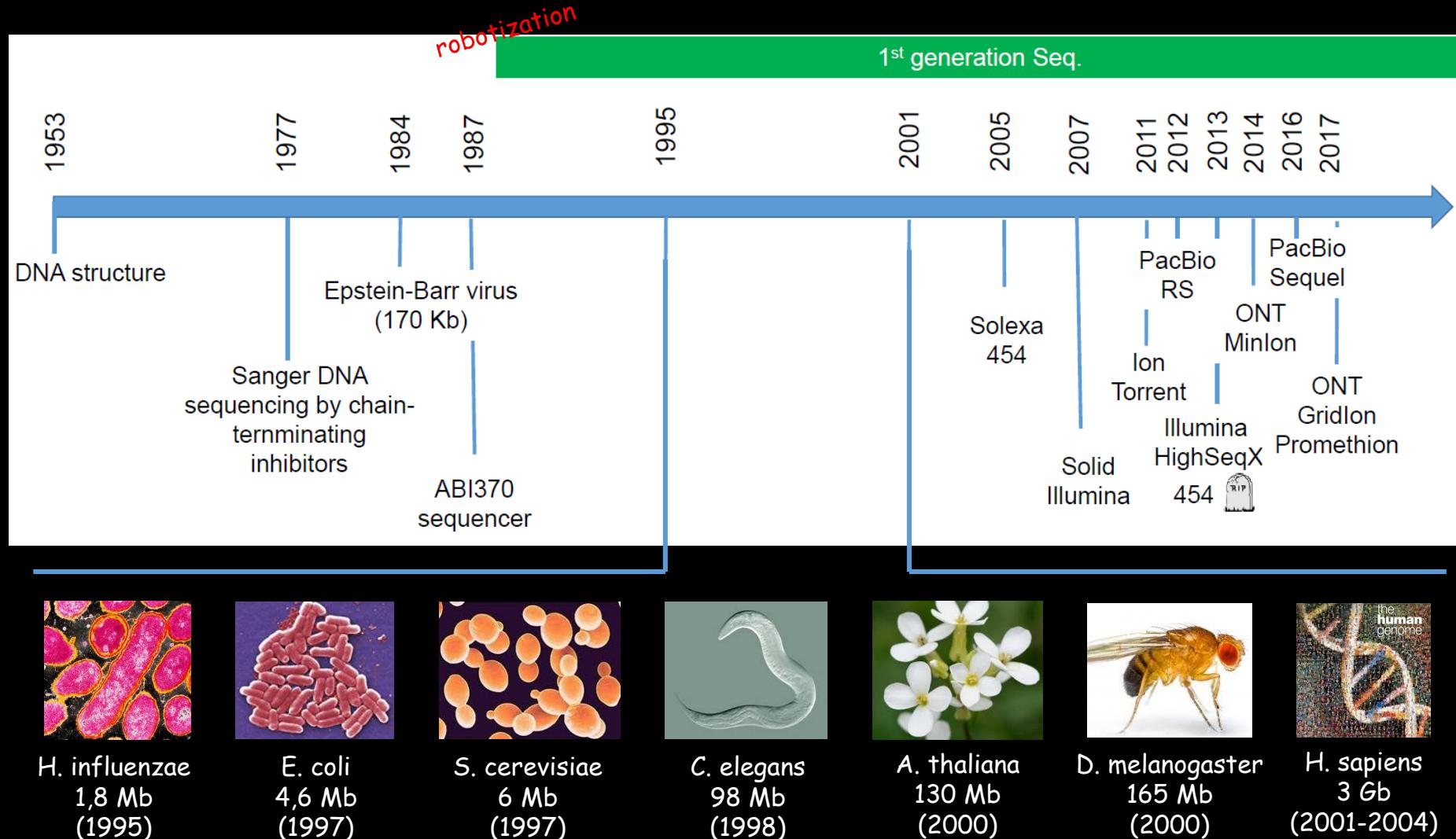
**Translational control**  
Initiation, elongation, termination  
Codon usage

**Post-translational control**  
Targeting  
Cleavage  
Modifications (glycosylation, acetylation...)

---

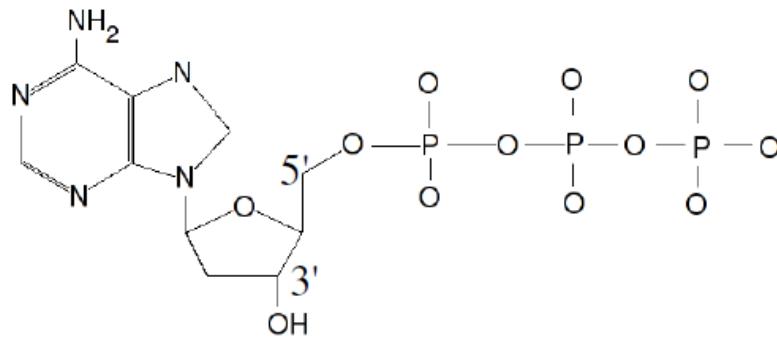
# Sequencing and Sequences

# History

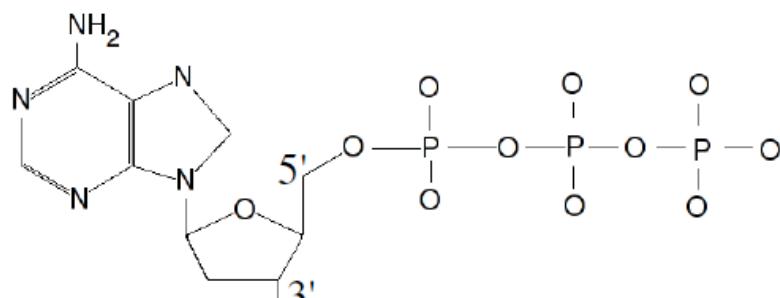


# Sanger DNA sequencing

Sanger sequencing is based on the selective incorporation of chain-terminating **dideoxynucleotides** by DNA polymerase during in vitro DNA replication



dATP

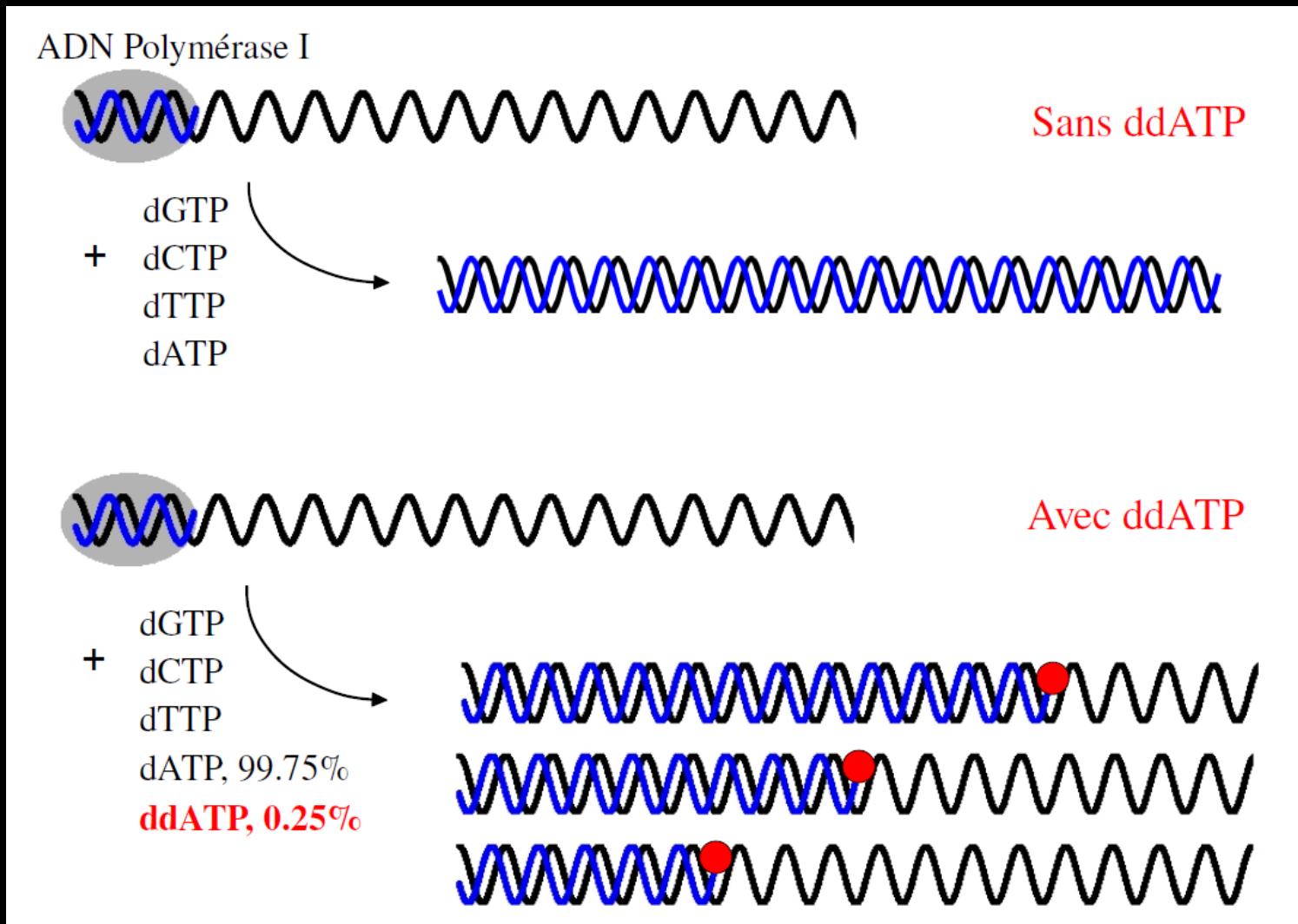


ddATP

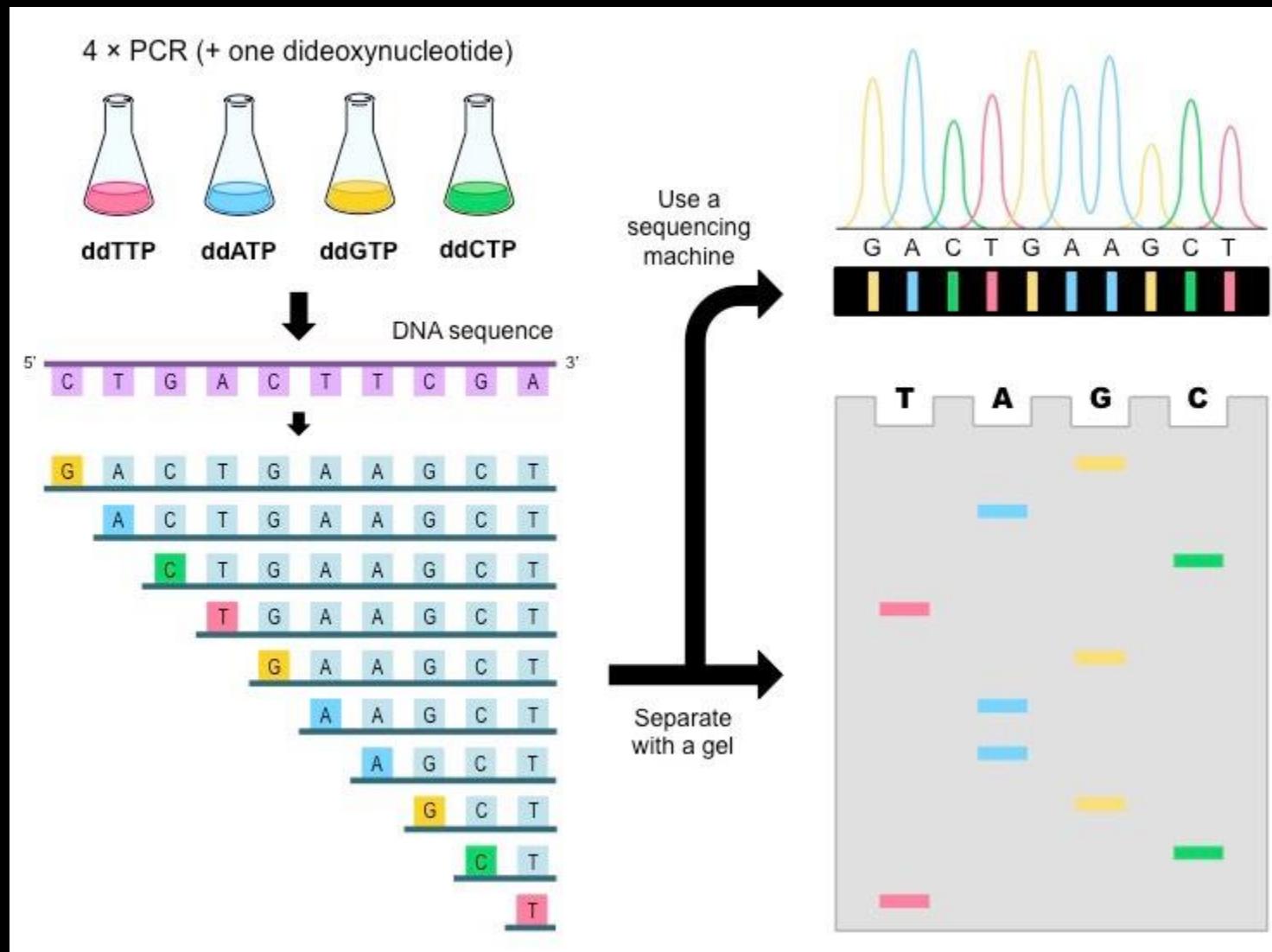
fluorophore



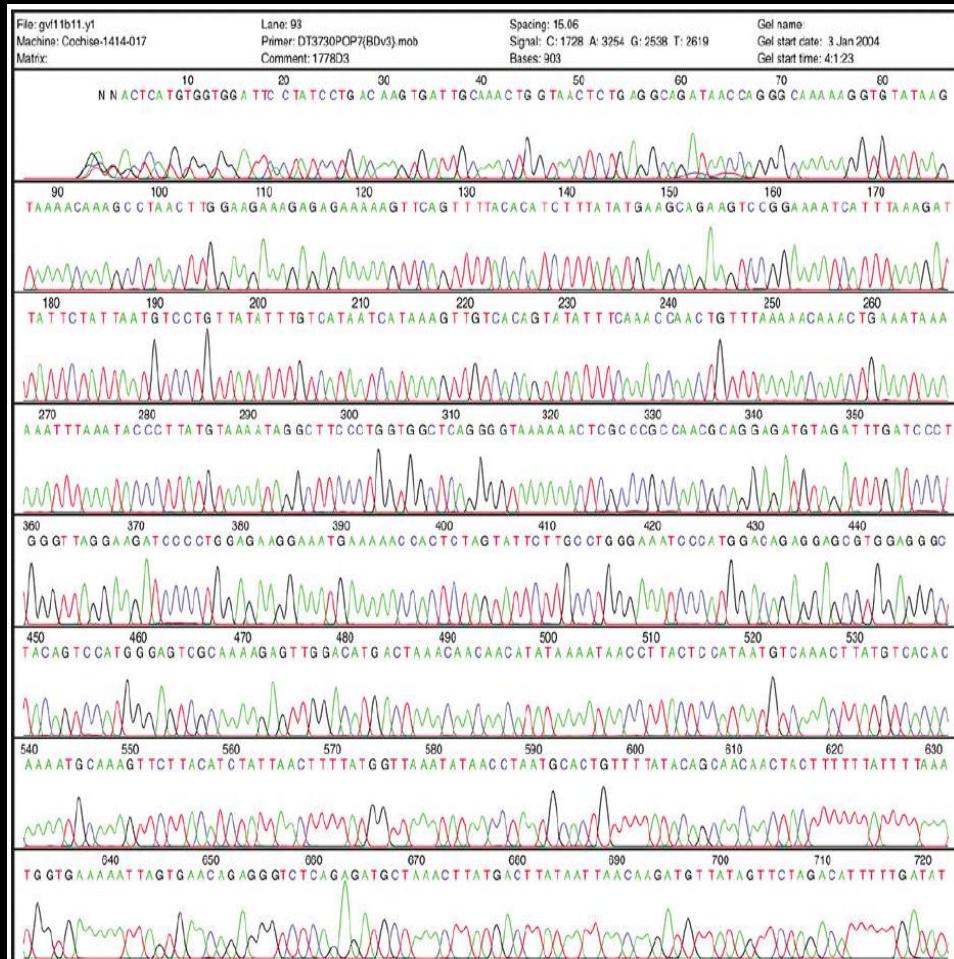
# Sanger DNA sequencing



# Sanger DNA sequencing

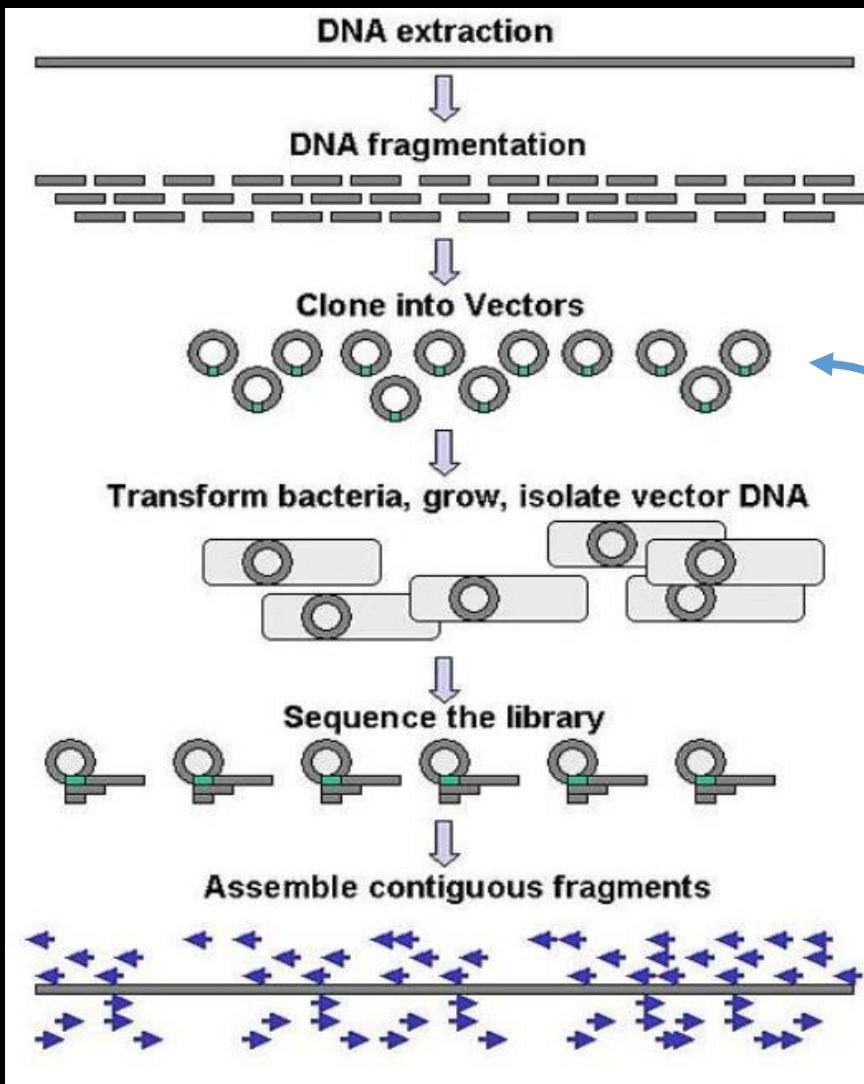


# Sanger DNA sequencing

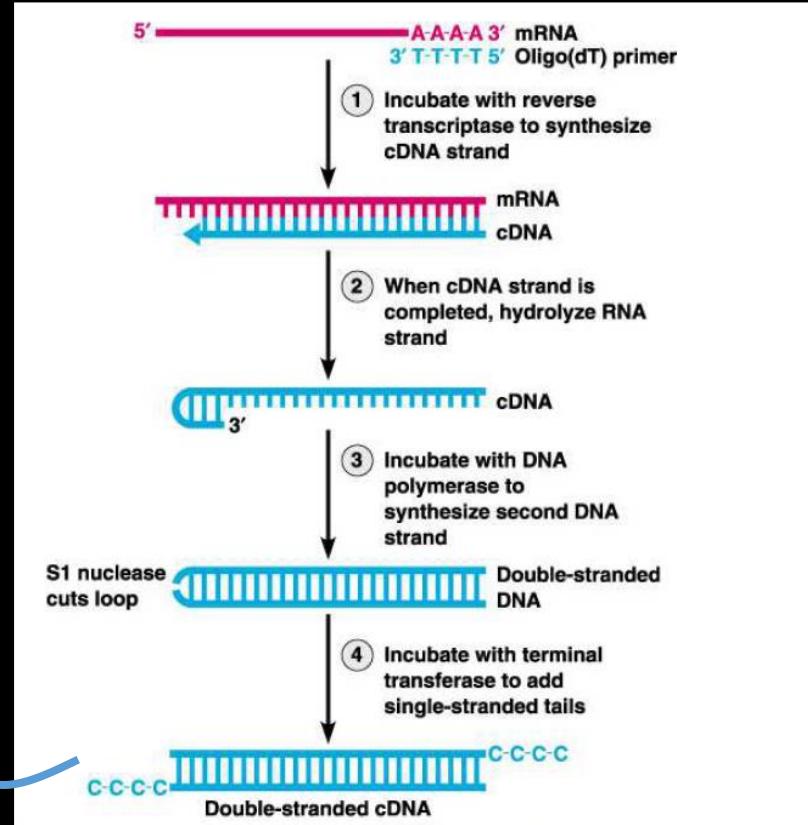


# A limiting cloning step

DNA

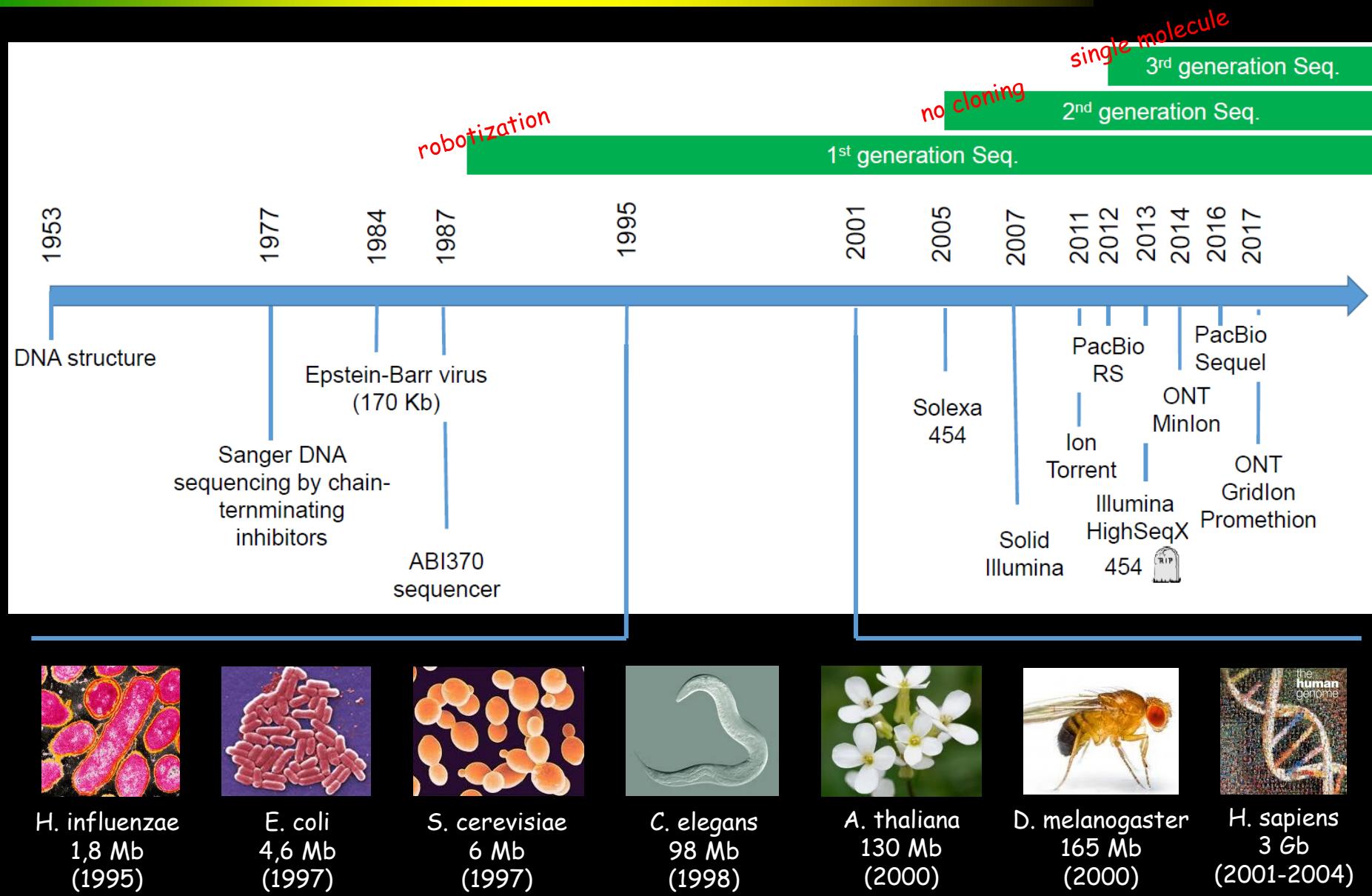


RNA

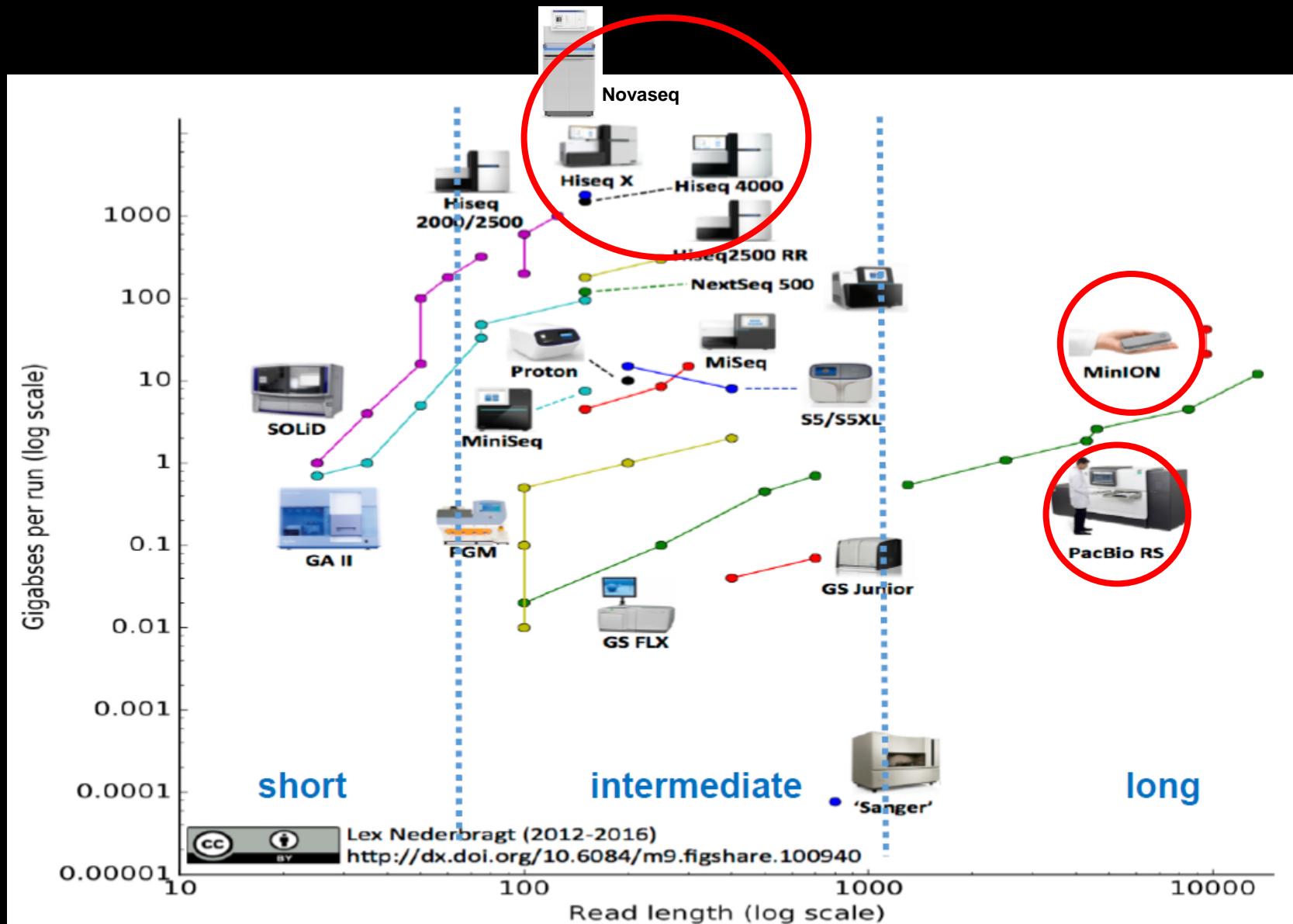


Reverse transcription

# History

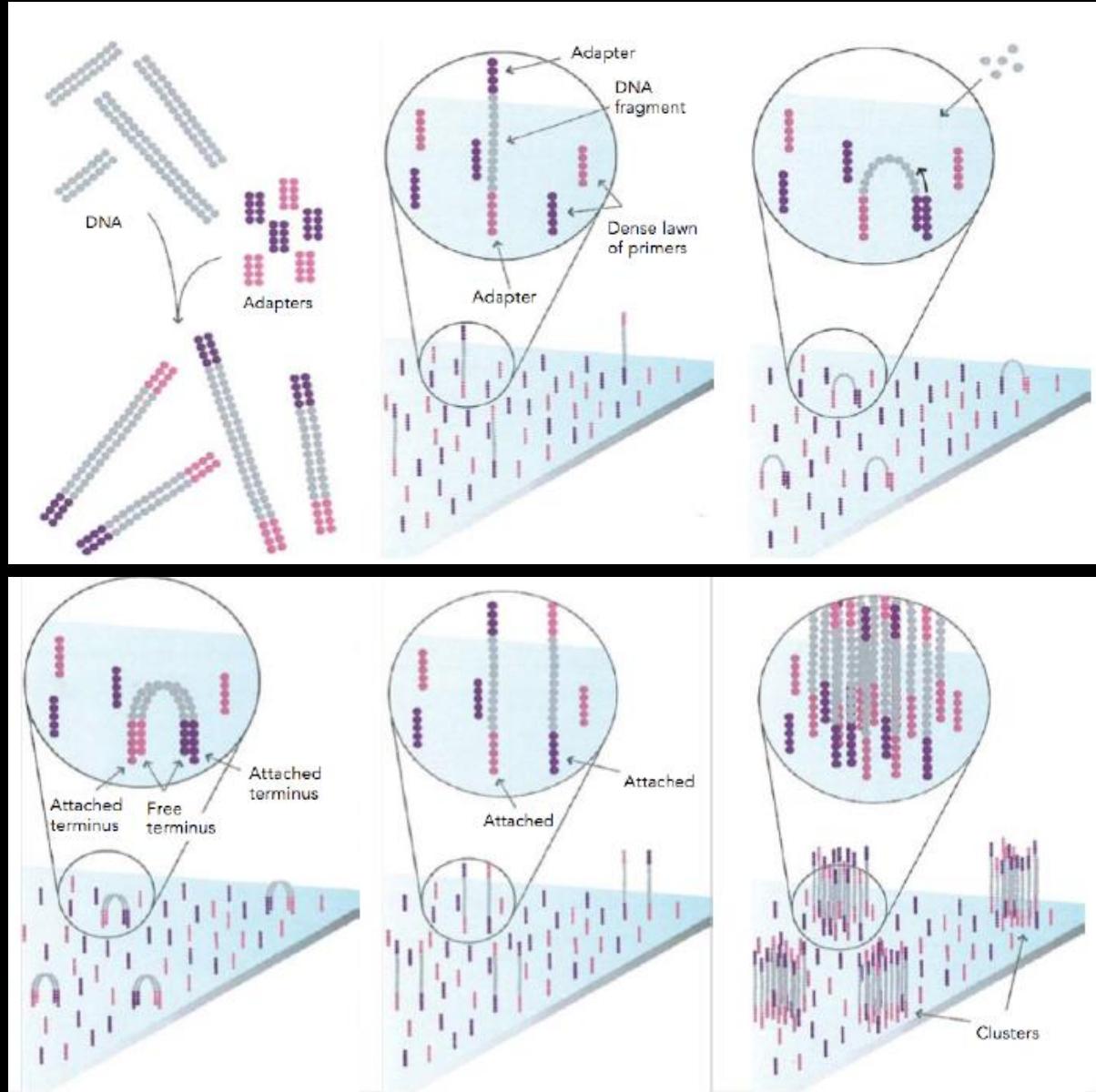


# Next generation sequencing



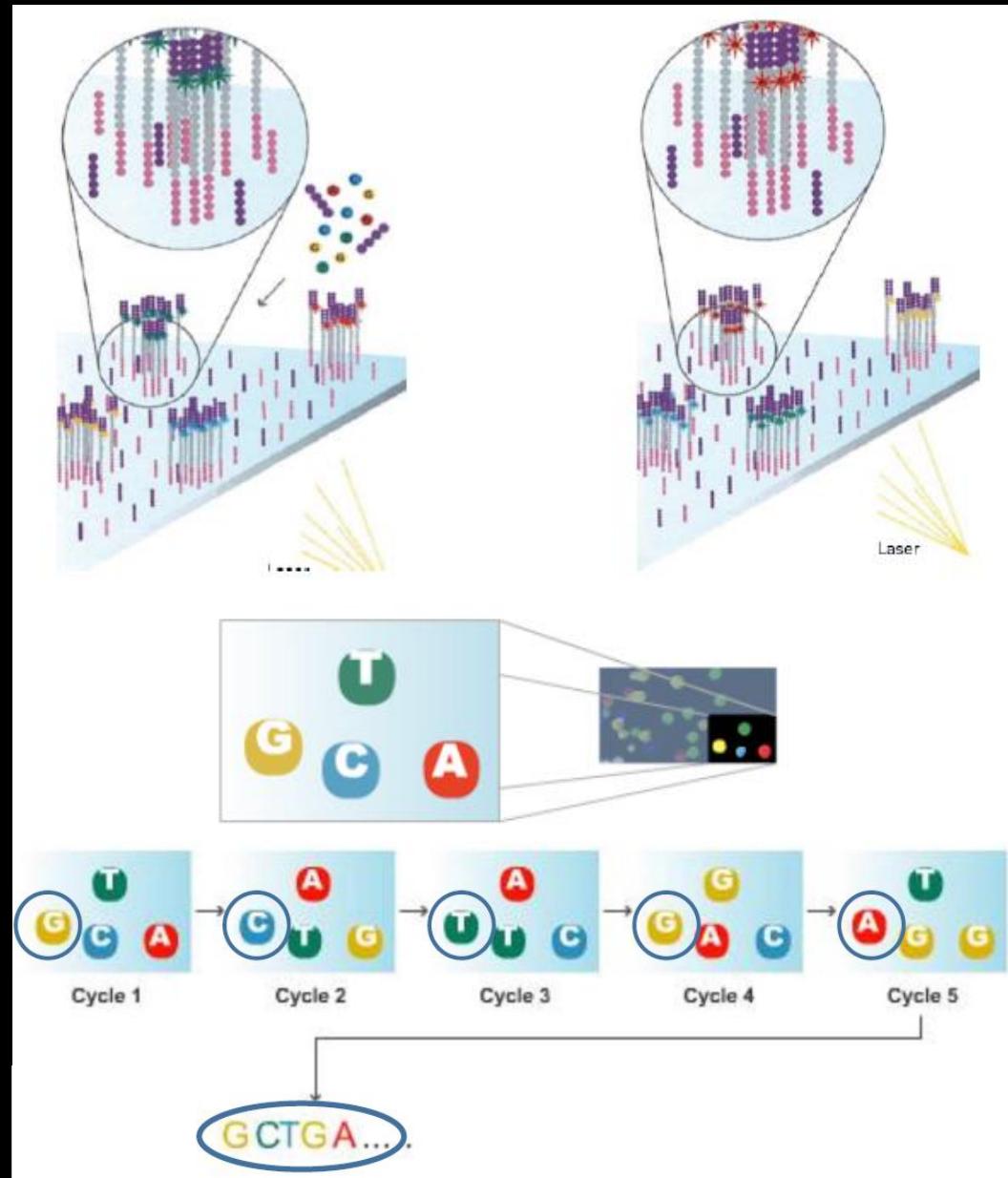
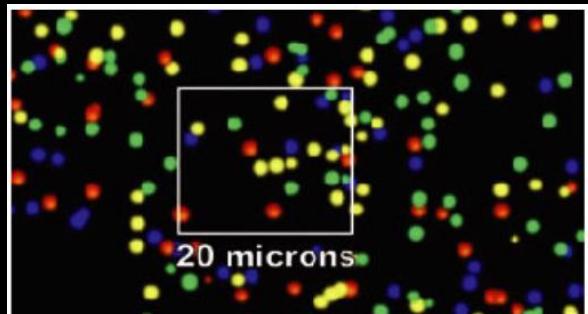
# Illumina technology (preparation)

- DNA library (fractioning, sizing)
- Adding specific adaptors
- DNA denaturation
- Random fixation on the flow cell
- PCR bridge
- Creation of clusters of amplified fragments



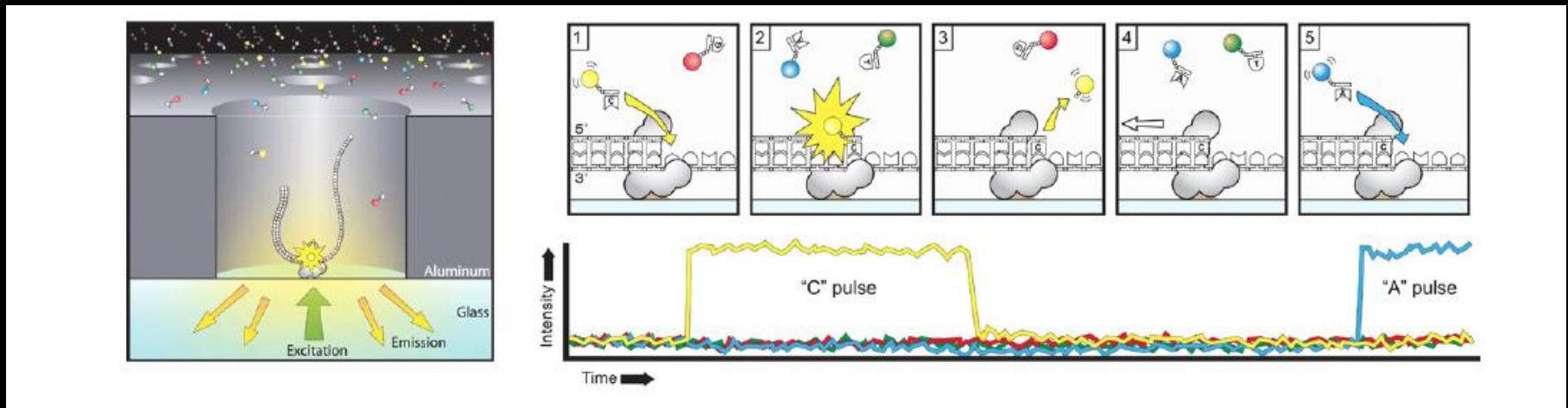
# Illumina technology (sequencing)

- Sequencing cycle with
  - 4 reversible terminators
  - labeled with fluorophore
  - primers
  - DNA-polymerase
- Laser excitation
- Image capture to identify the incorporated base for each cluster
- Next cycle after removing the blocked 3' terminus and the fluorophore from each added base

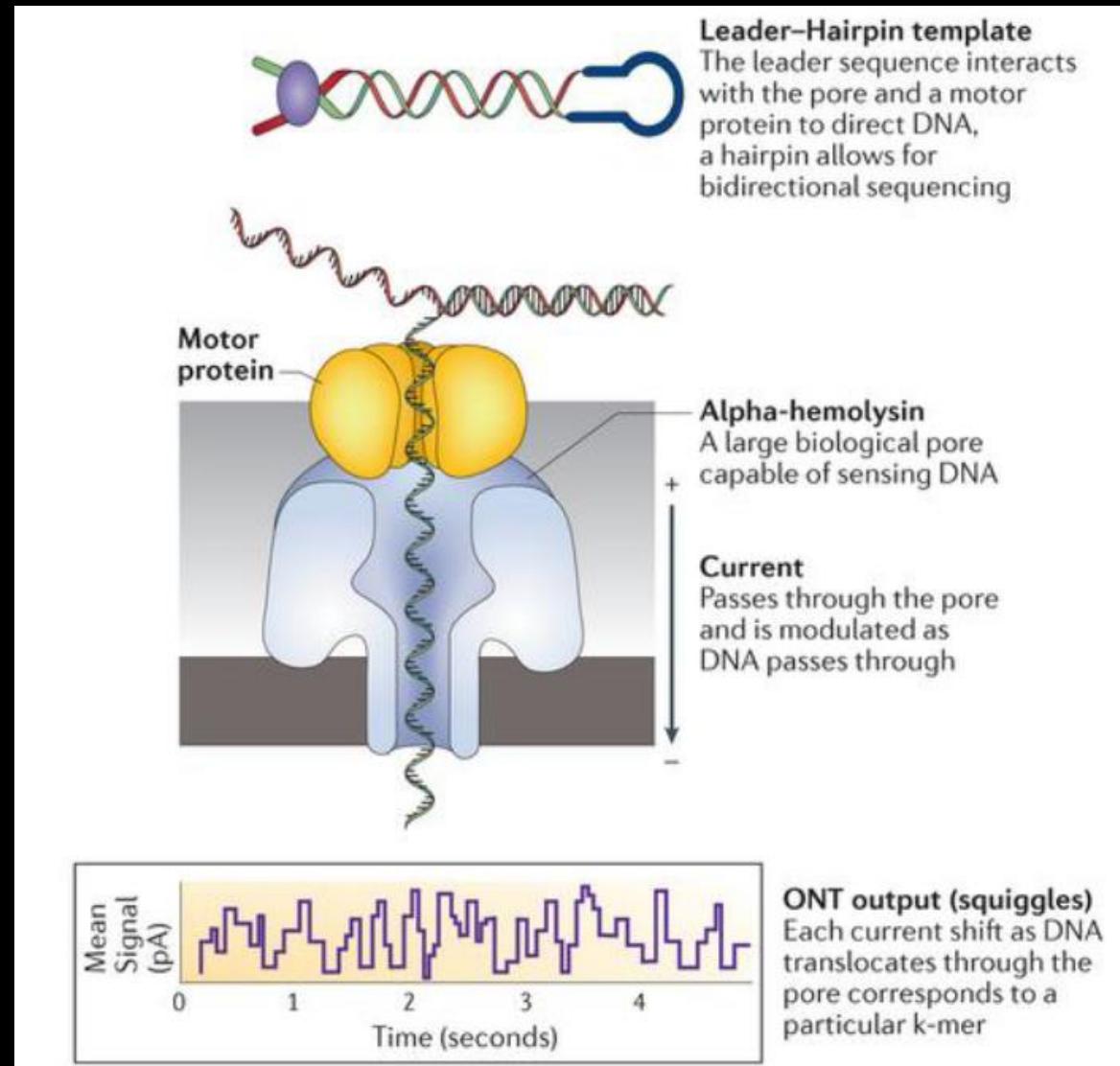


# PacBio technology

Single Molecule, Real-Time (SMRT) sequencing / Long-read sequencing  
DNA-polymerase is fixed in a micro-well ( $10^{-21}$  liters)  
Labeled nucleotides incorporation is detected in real-time

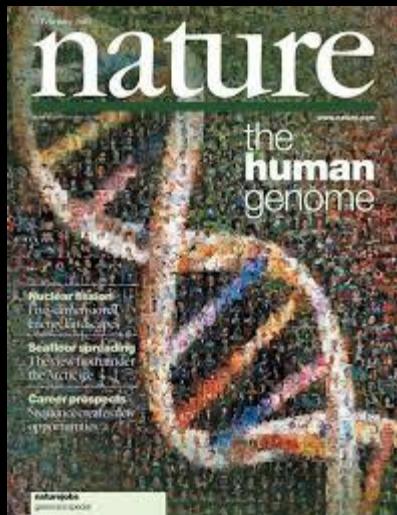


# Oxford Nanopore technology



# Some numbers...

Technology	Read length (b)	Reads per run	Bases per run (Gb)	Error rate
ABI Sanger	800	96	0,0000768	< 1%
Illumina	150	$600.10^6$ to $6.10^9$	7 to 2000	< 1%
PacBio	13 500	660 000	12 000	10-15%
Nanopore	9 545	$4,4.10^6$	42	5-15%

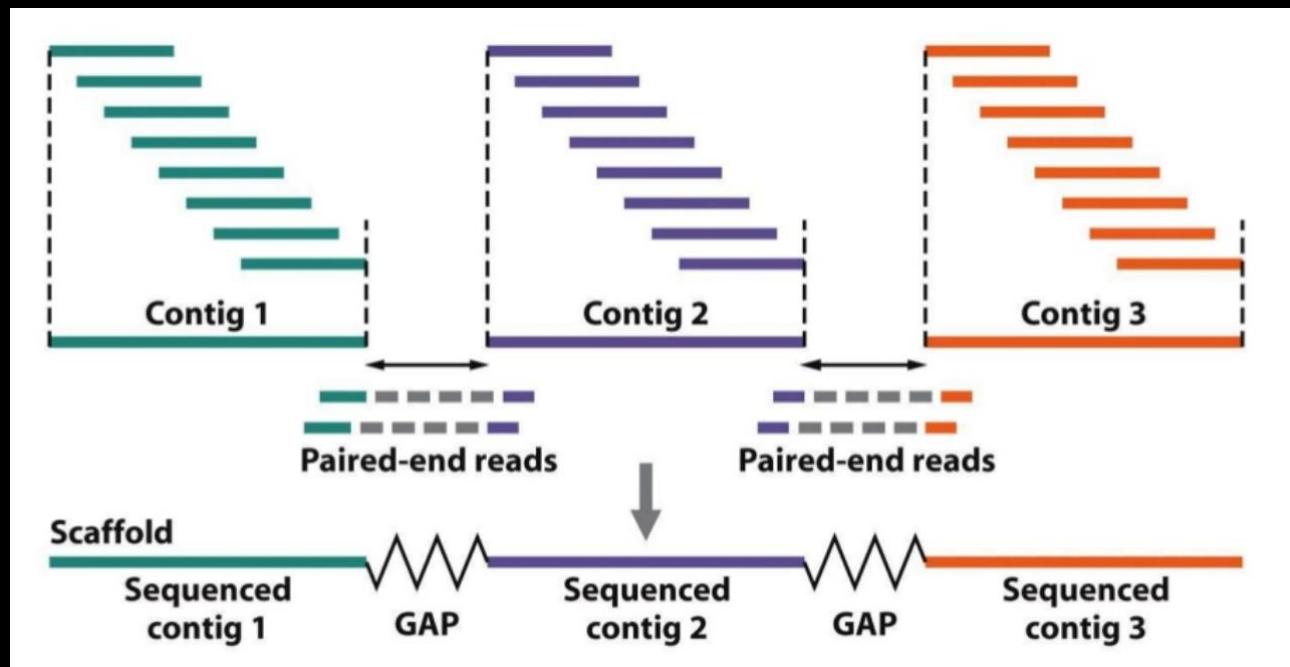


Human Genome Project

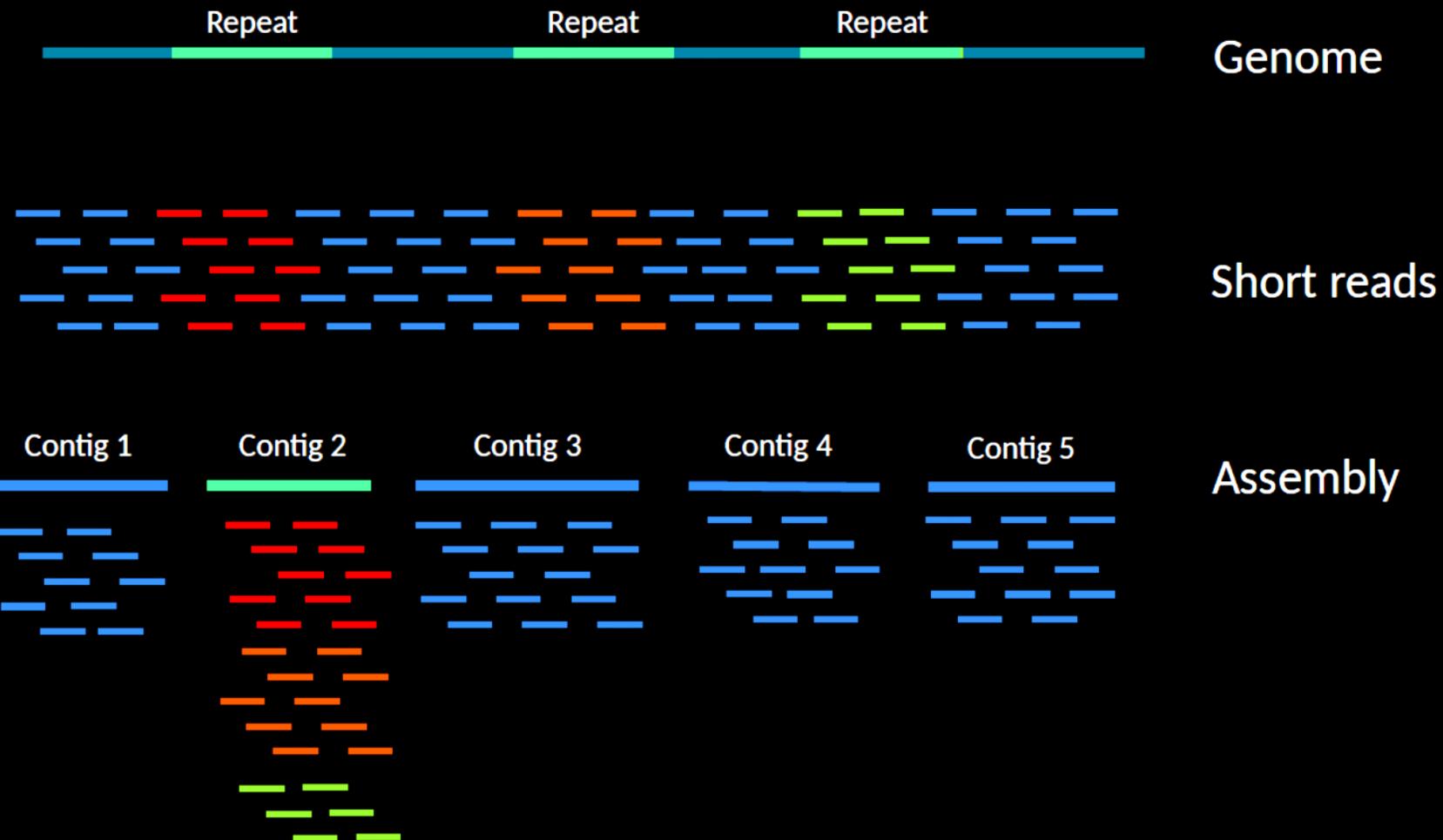
13 years (1990-2003)  
\$ 3,000,000,000

Illumina XTen technology (2014)  
1 week  
\$ 1,000

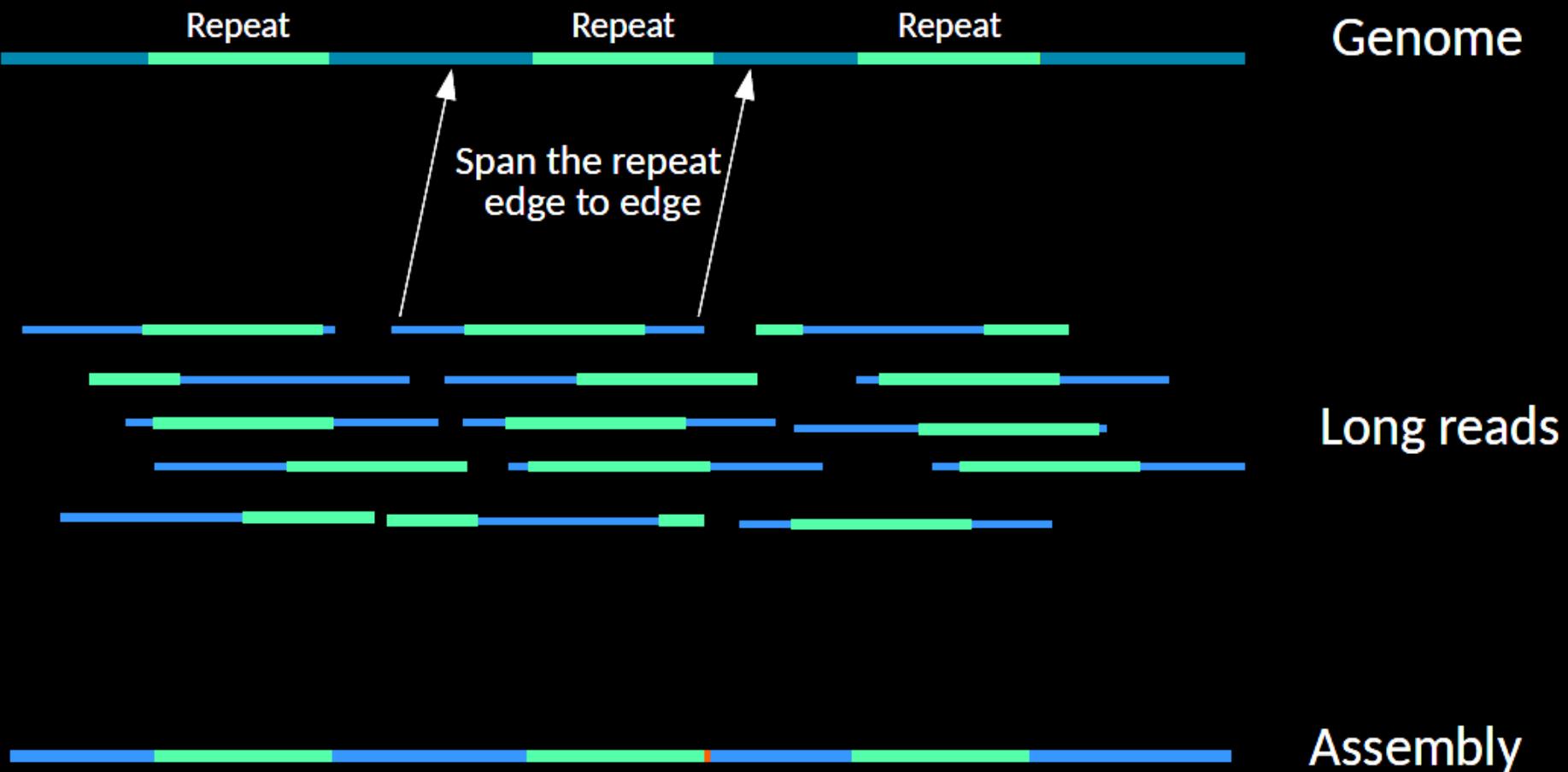
# Genome assembly : the challenge



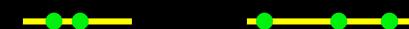
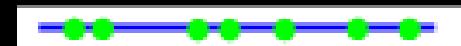
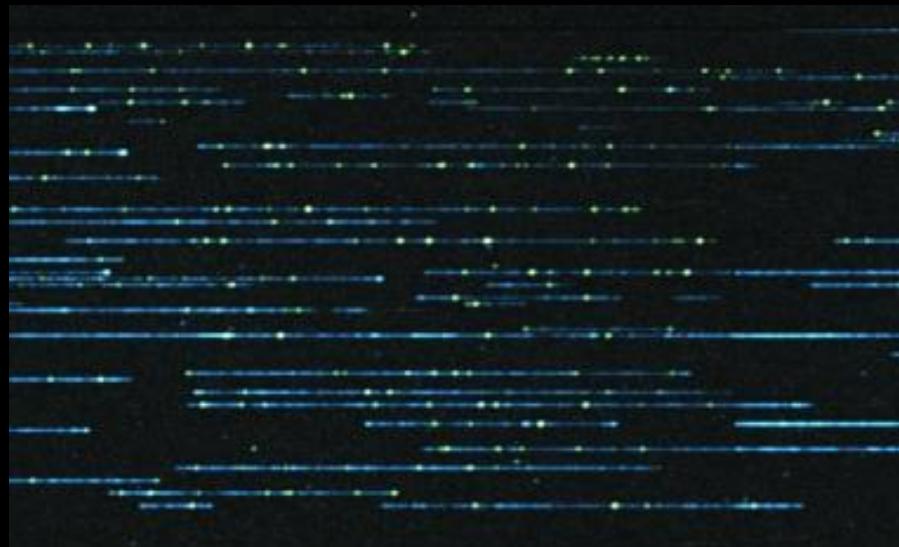
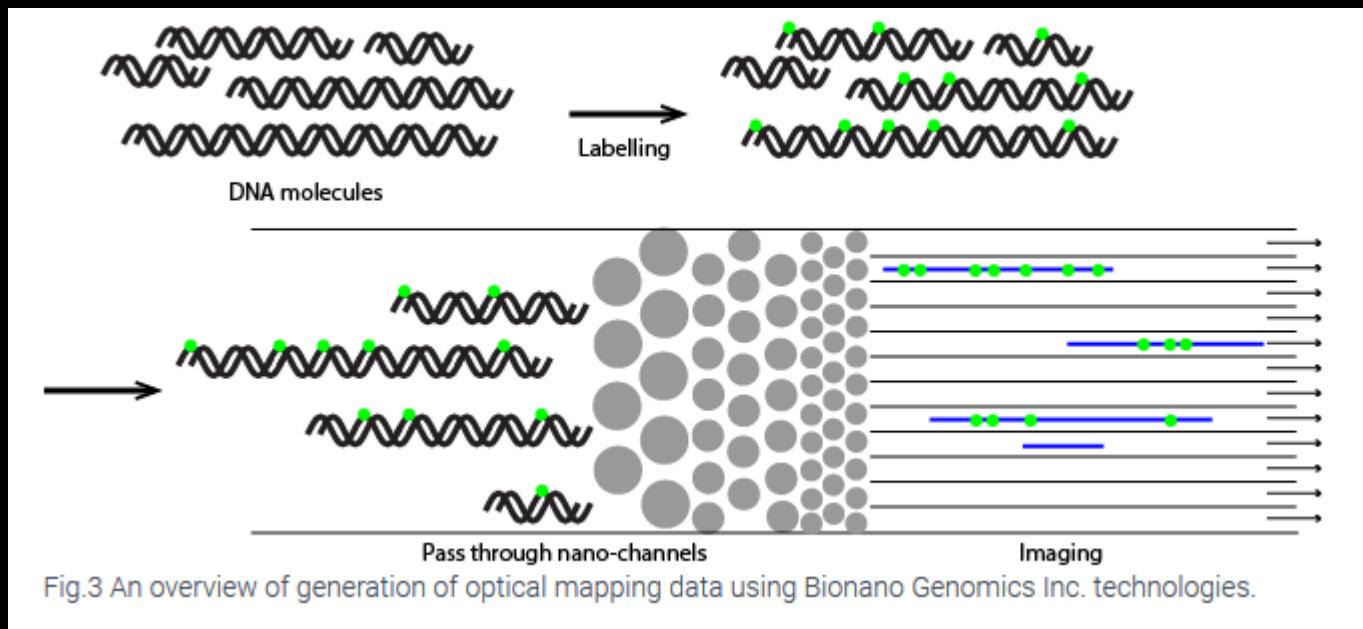
# Genome assembly : advantage of long reads



# Genome assembly : advantage of long reads

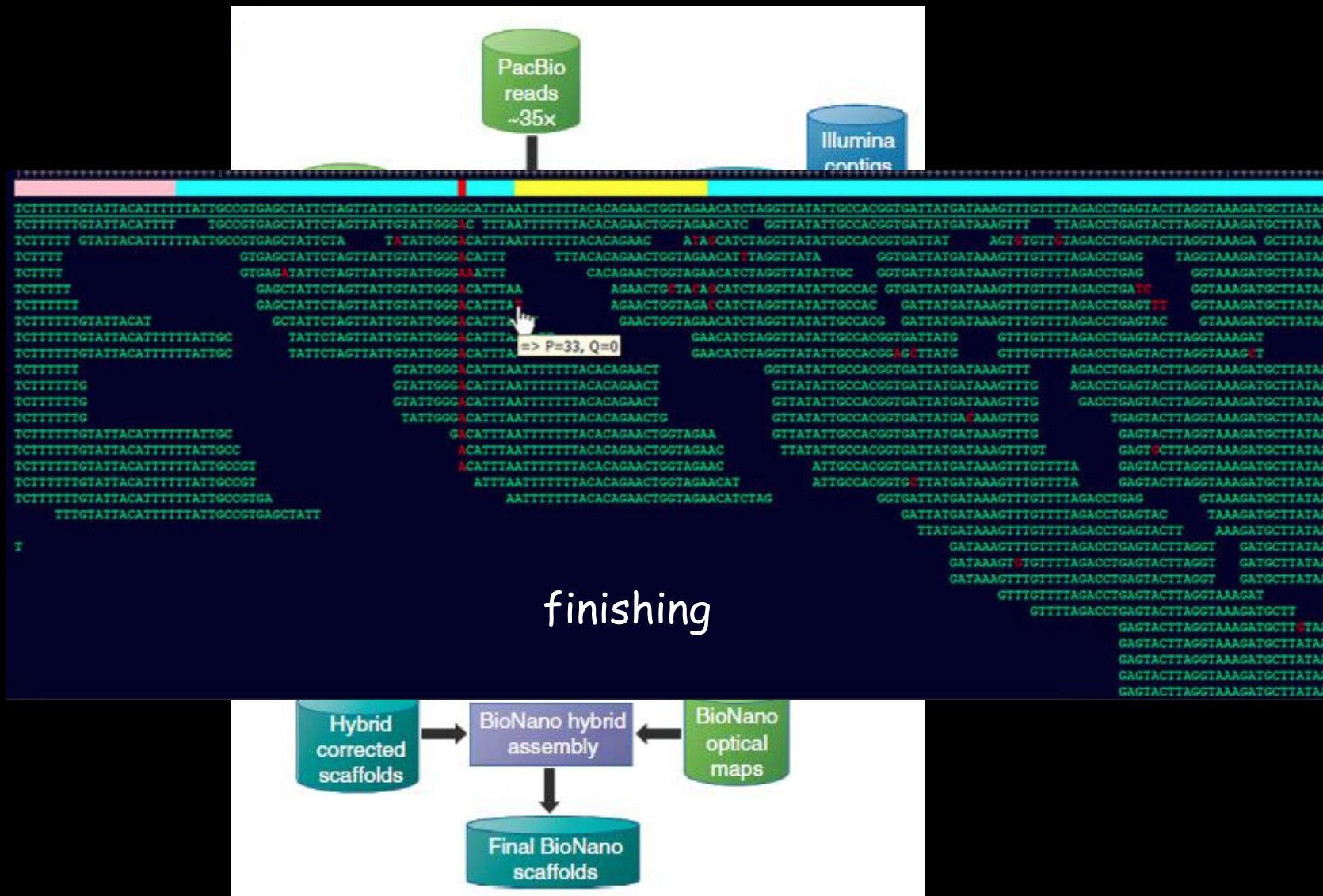


# Genome assembly : optical mapping



scaffolding

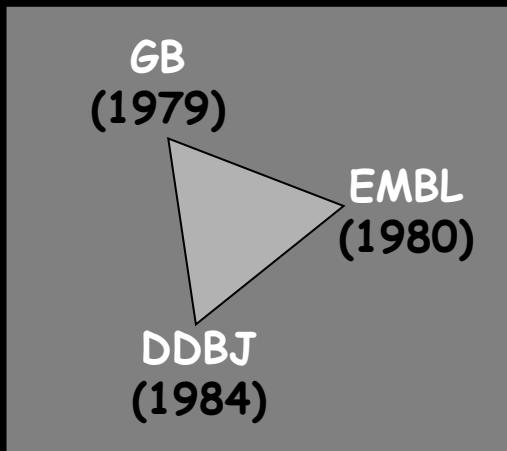
# Genome assembly : Hybrid assembly strategy



# Bank of sequences

---

## International repositories



synchronized

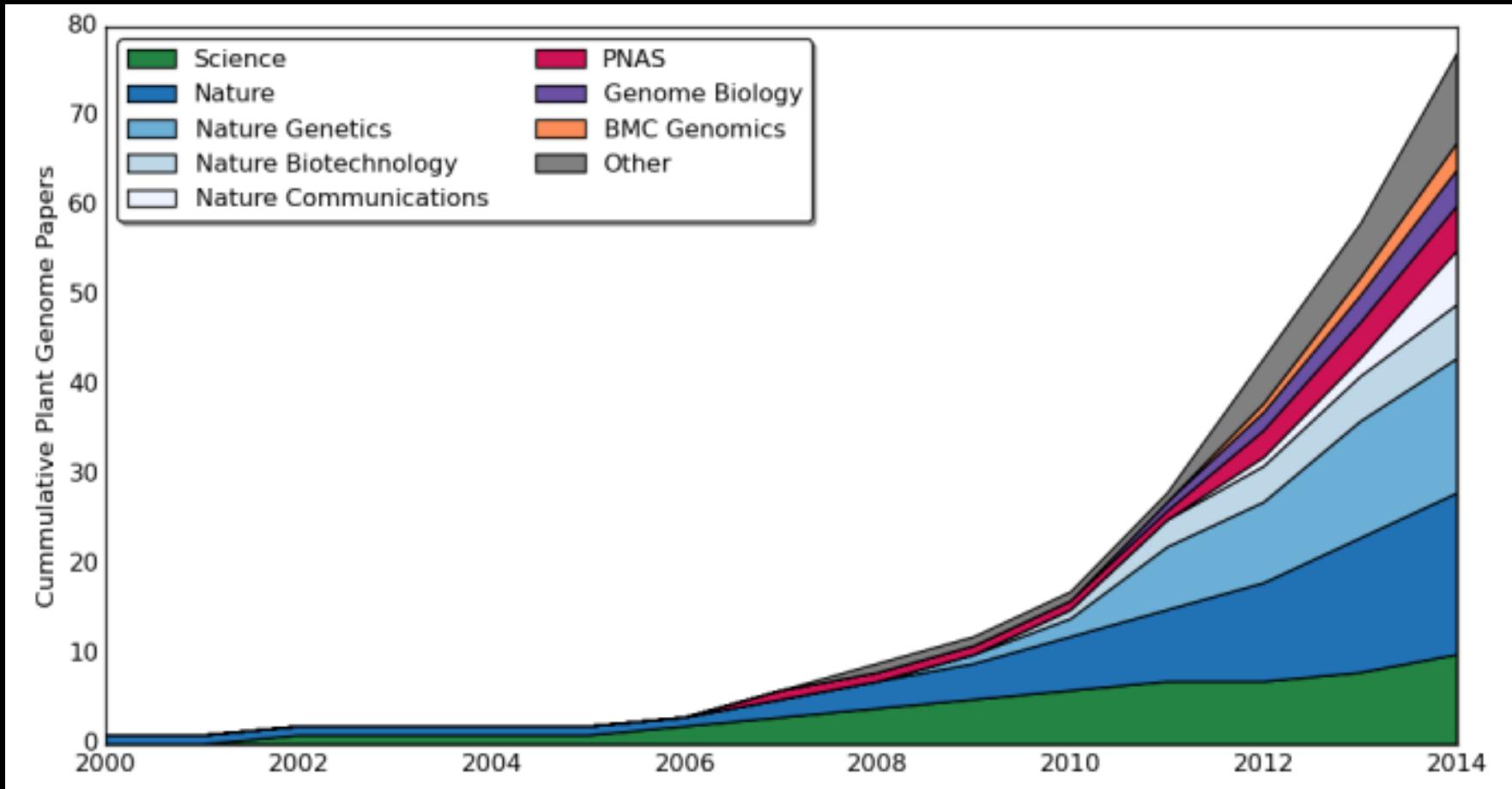
3 different formats

- GENBANK @ NCBI (National Center for Biotechnology Information)
- EMBL (European Molecular Biology Laboratory) @ EBI (European Bioinformatics Institute)
- DDBJ (DNA DataBase of Japan) @ DDBJ center

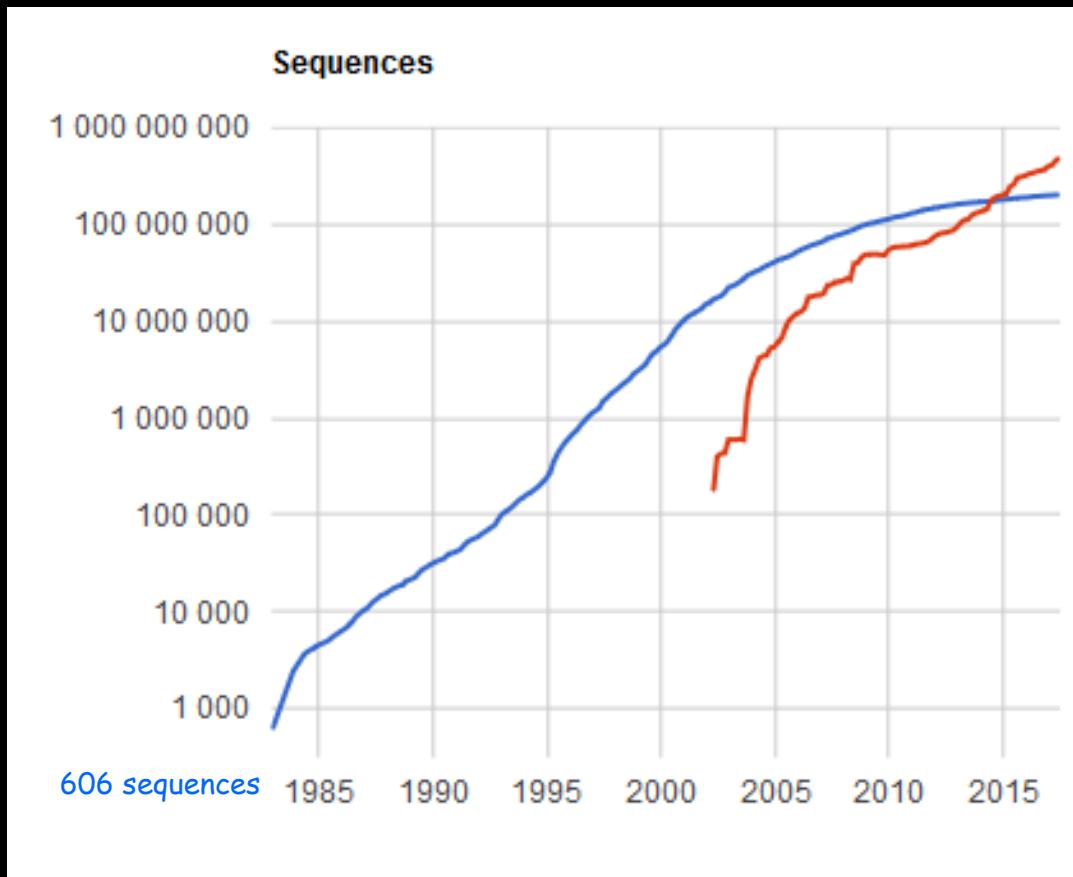
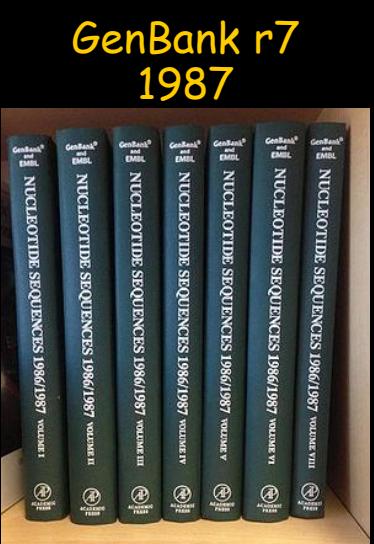
# High throughput sequencing

## Genomics sizes up

NATURE|Vol 451|17 January 2008



# High throughput sequencing



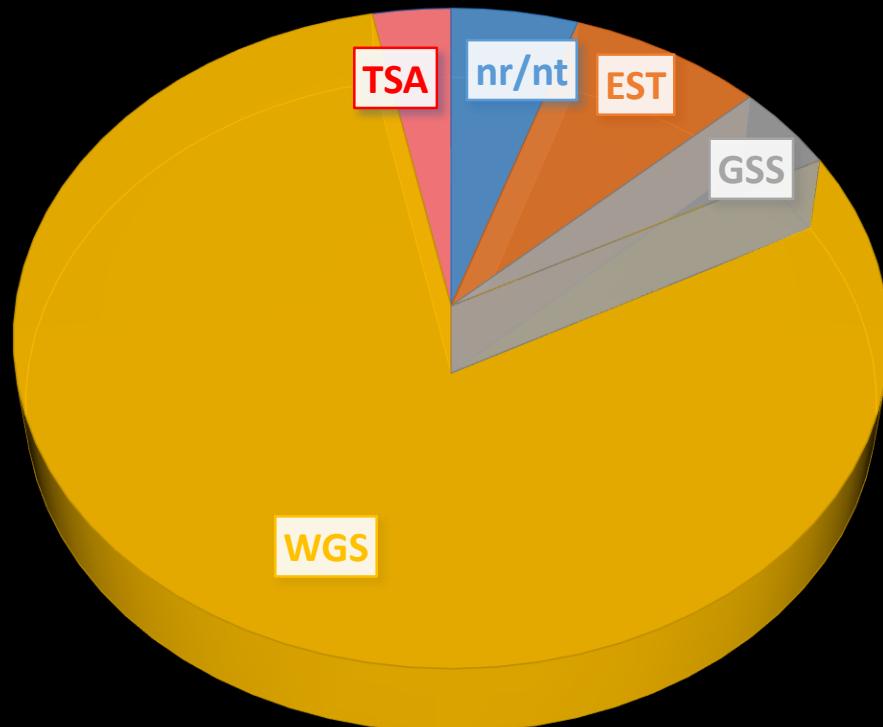
WGS  
 $3656.10^9$  nt

GENBANK  
 $285.10^9$  nt

release occurs every 2 months  
doubling every 15 months

# GenBank divisions

Huge diversity of sequences with different **origins** (species, mutants), **lengths** (few bases to whole chromosomes), **qualities** (clones, assembly, single-pass, gold finishing), **types** (transcripts, genomic sequences, metagenomes, synthetics)...



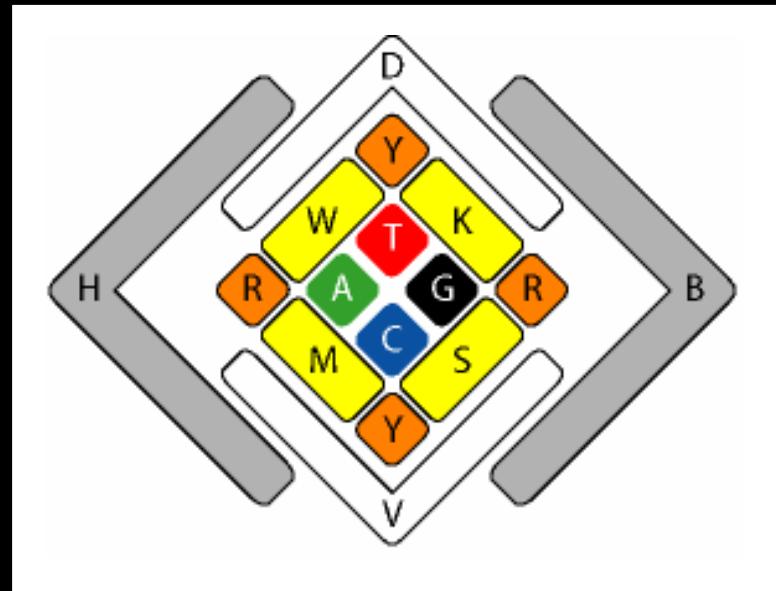
millions of sequences

nr/nt	50
EST	77
GSS	41
TSA	31
WGS	773

# GenBank notes

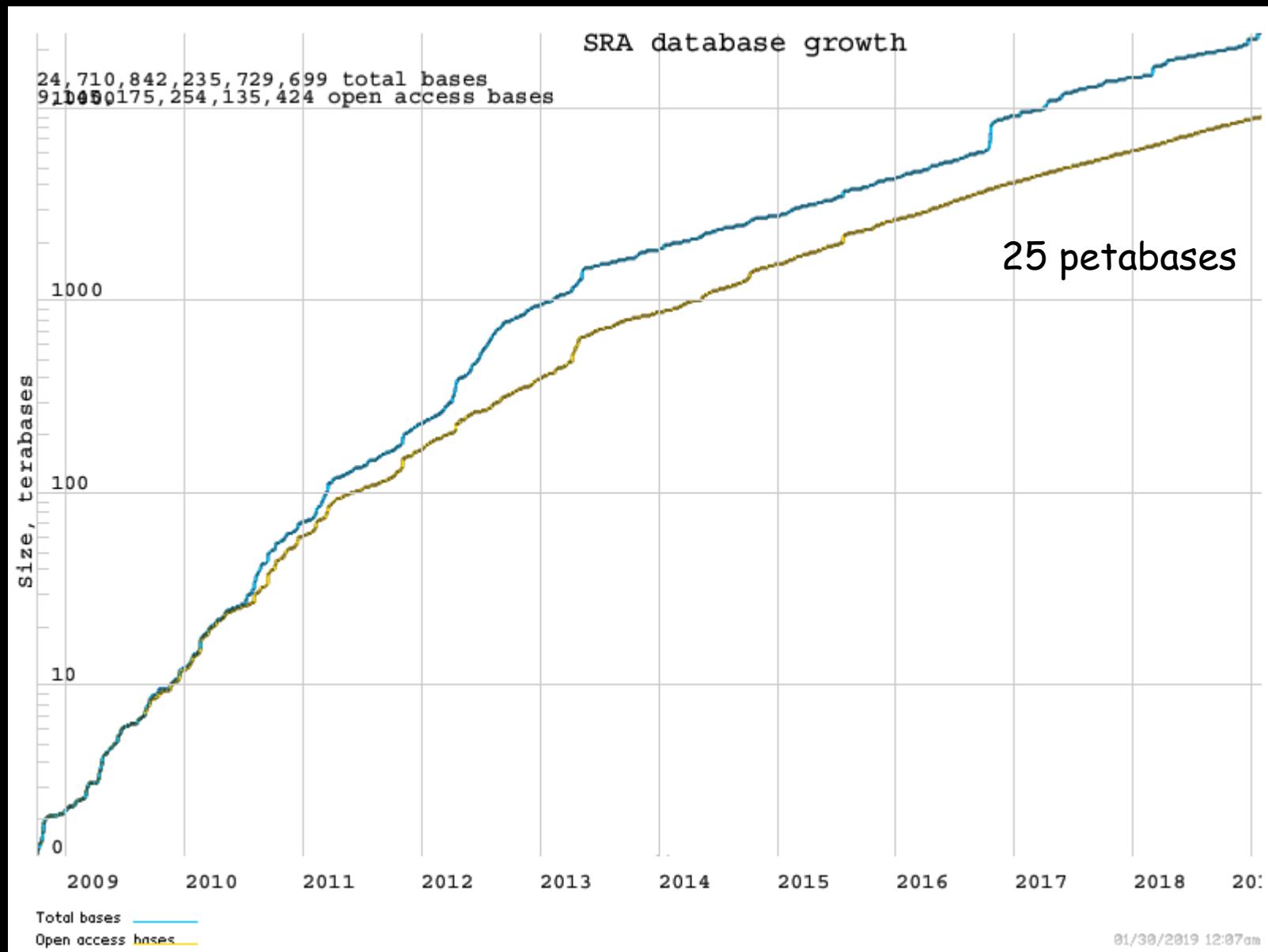
---

All sequences (RNA or DNA) contains A, T, G, C and degenerated IUB code but never U.



The strand available in GenBank is named '+' or 'direct' without consideration regarding the coding strand.

# Sequence Read Archive (SRA division)



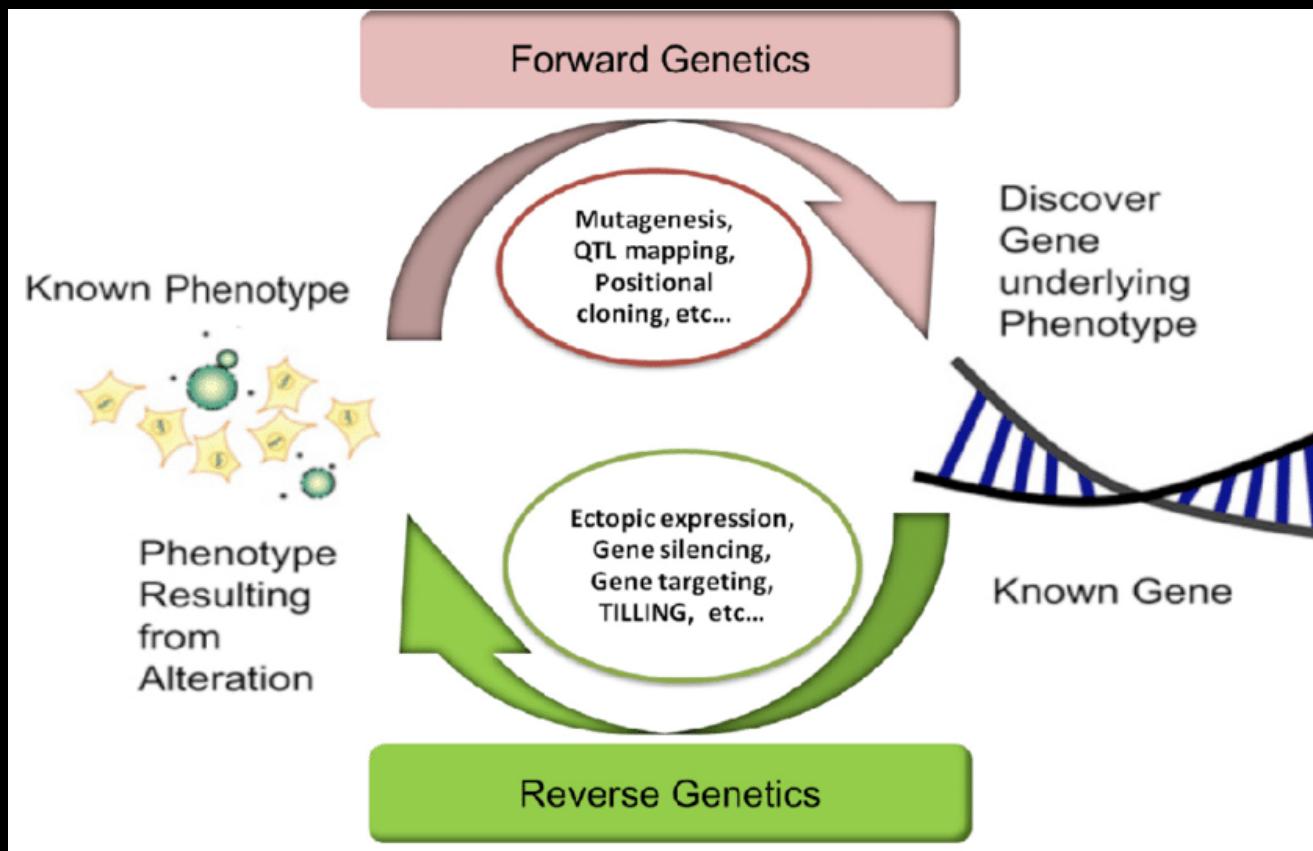
# From genome to genes... from genes to functions



# Genomics : from forward to reverse genetics

Forward 'classical' genetics seeks to find the genetic base of a phenotype or trait.

Reverse genetics seeks to find what phenotype arise a result of particular genetic sequence.



---

# Genome Annotation

# Annotation

Sequences  $\longleftrightarrow$  Biological information

## Structural Annotation

- gene localization ?
- gene structure ?
- coding regions ?
- regulatory sites ?

## Functional Annotation

- biochemical function ?
- biological functions ?
- regulation and interactions ?
- expression ?



Experimental work

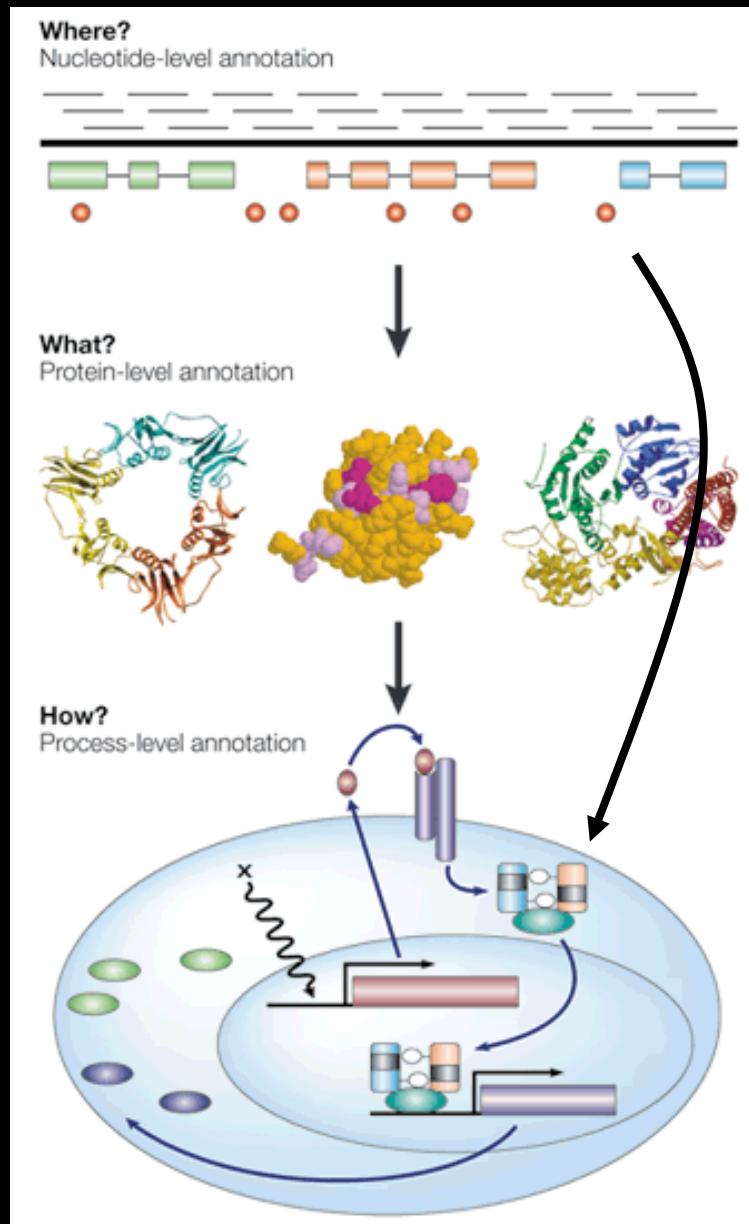


*in silico* predictions

# Annotation

Structural Annotation

Functional Annotation



biochemical  
function

biological  
function

Relational  
Annotation

Systems  
Biology

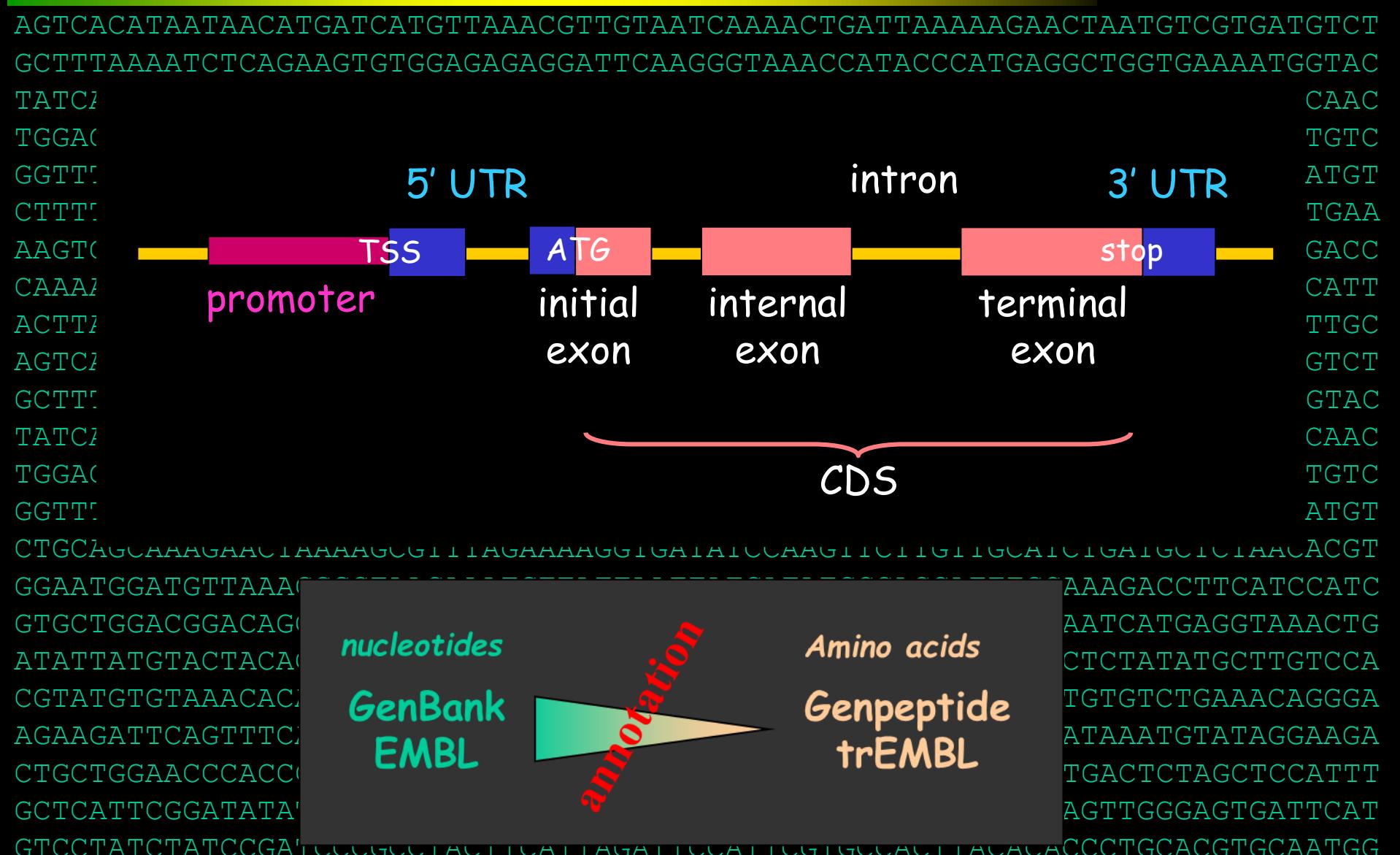


---

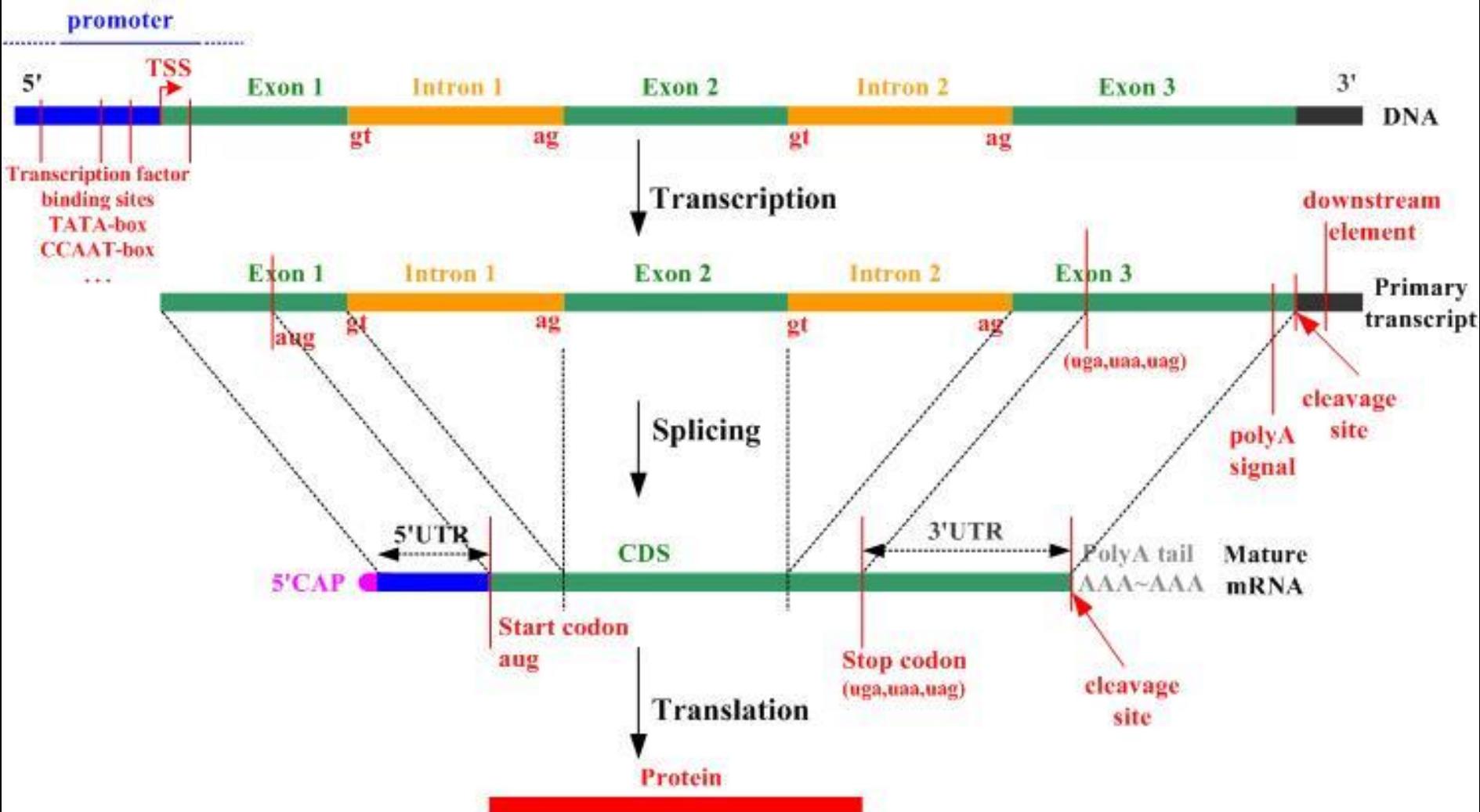
# Genome Annotation

# Structure

# Structural annotation: goals

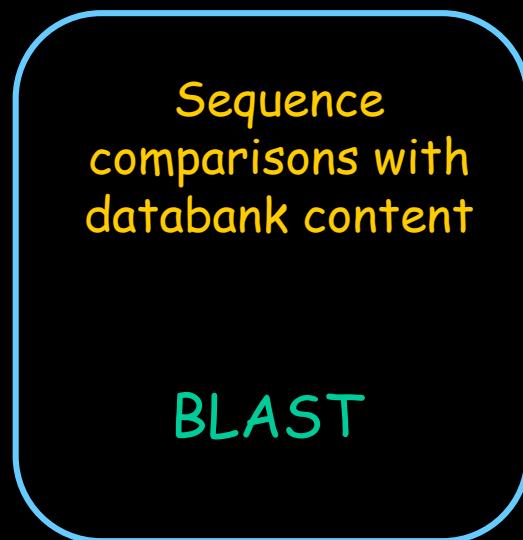


# As a reminder...

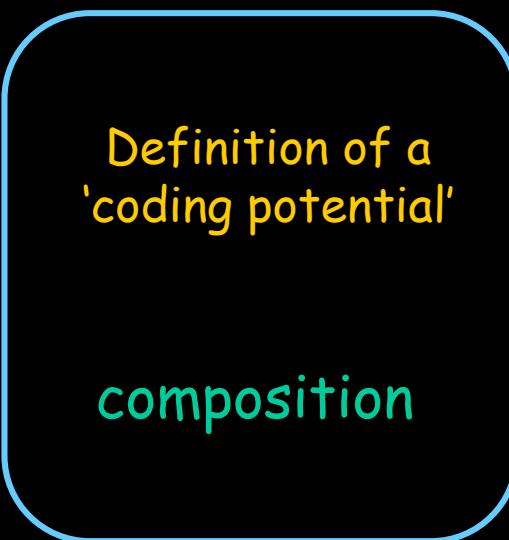


# Prediction methods

by similarity



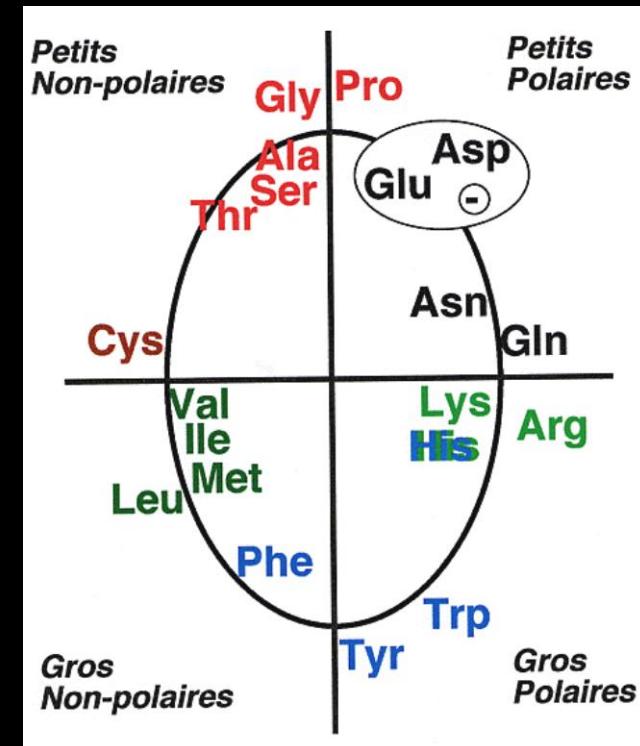
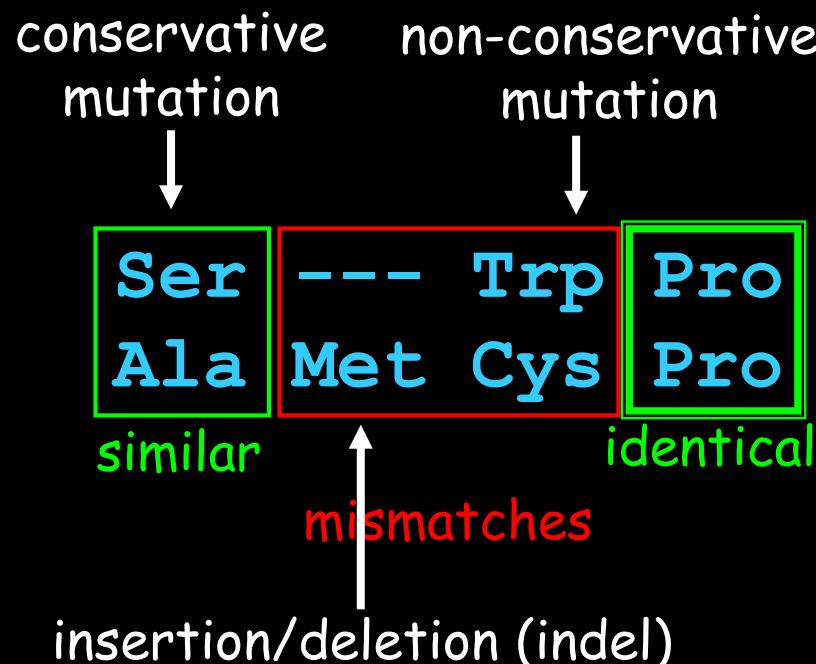
*ab initio*



Learning step is necessary

# Quantification of similarities

Representation of a DAYHOFF matrix of mutations (substitutions) by a plan projection.



# Quantification of similarities

PAM (1978)

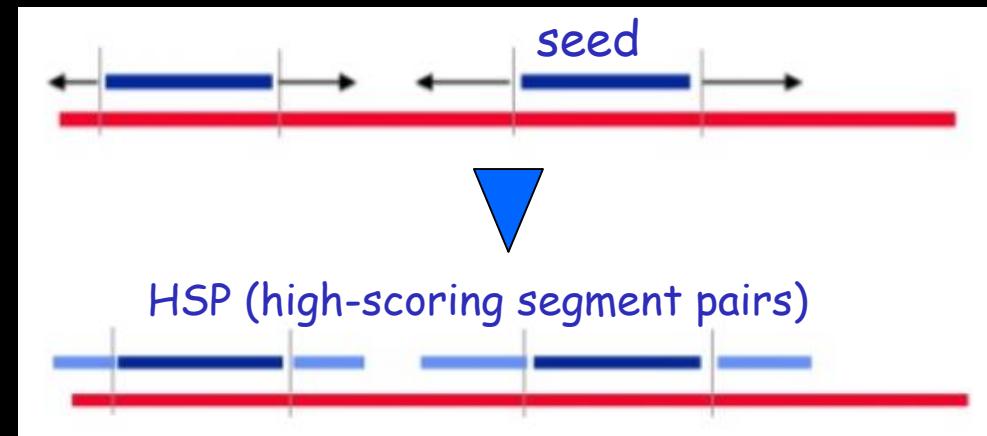
BLOSUM (1992)

BLOSUM62

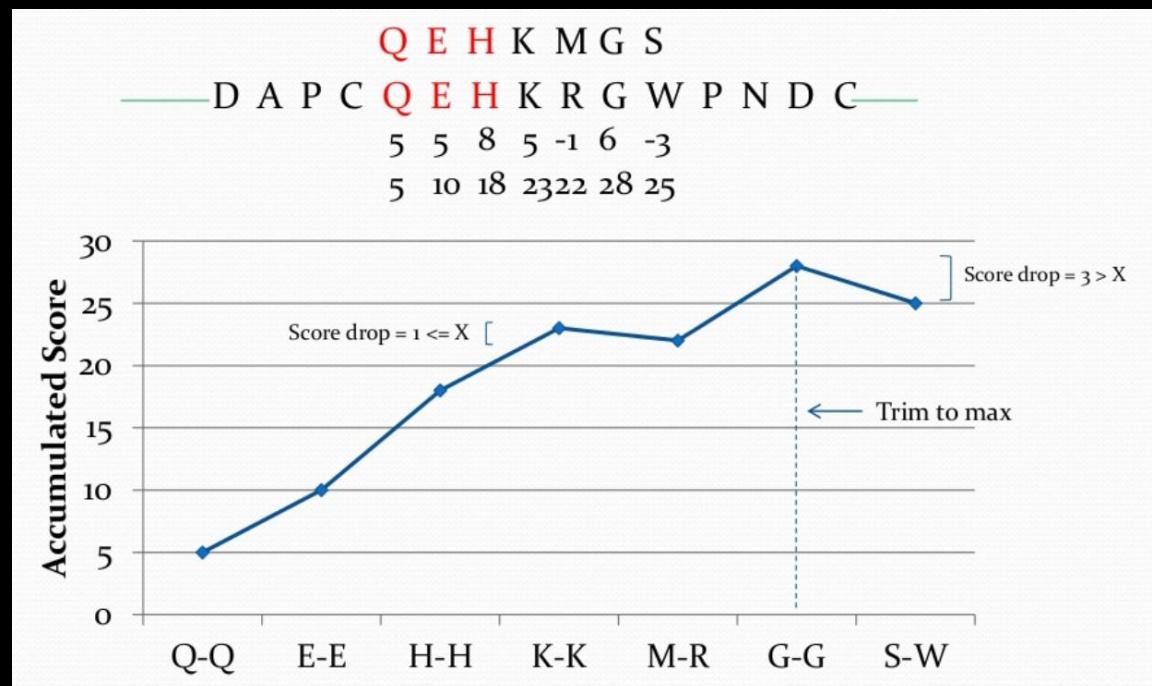
	C	G	A	T	S	N	D	E	Q	K	R	V	I	L	M	W	F	Y	H	P	
C	9	-3	0	-1	-1	-3	-3	-4	-3	-3	-3	-1	-1	-1	-1	-2	-2	-2	-3	-3	
G	-3	6	0	-2	0	0	-1	-2	-2	-2	-2	-3	-4	-4	-3	-2	-3	-3	-2	-2	
A	0	0	4	0	1	-2	-2	-1	-1	-1	-1	0	-1	-1	-1	-3	-2	-2	-2	-1	
T	-1	-2	0	5	1	0	-1	-1	-1	-1	-1	0	-1	-1	-1	-2	-2	-2	-2	-1	
S	-1	0	1	1	4	1	0	0	0	0	0	-1	-2	-2	-2	-1	-3	-2	-2	-1	-1
N	-3	0	-2	0	1	6	1	0	0	0	0	-3	-3	-3	-2	-4	-3	-2	1	-2	
D	-3	-1	-2	-1	0	1	6	2	0	-1	-2	-3	-3	-4	-3	-4	-3	-3	-1	-1	
E	-4	-2	-1	-1	0	0	2	5	2	1	0	-2	-3	-3	-2	-3	-3	-2	0	-1	
Q	-3	-2	-1	-1	0	0	0	2	5	1	1	-2	-3	-2	0	-2	-3	-1	0	-1	
K	-3	-2	-1	-1	0	0	-1	1	1	5	2	-2	-3	-2	-1	-3	-3	-2	-1	-1	
R	-3	-2	-1	-1	-1	0	-2	0	1	2	5	-3	-3	-2	-1	-3	-3	-2	0	-2	
V	-1	-3	0	0	-2	-3	-3	-2	-2	-2	-3	4	3	1	1	-3	-1	-1	-3	-2	
I	-1	-4	-1	-1	-2	-3	-3	-3	-3	-3	-3	3	4	2	1	-3	0	-1	-3	-3	
L	-1	-4	-1	-1	-2	-3	-4	-3	-2	-2	-2	1	2	4	2	-2	0	-1	-3	-3	
M	-1	-3	-1	-1	-1	-2	-3	-2	0	-1	-1	1	1	2	5	-1	0	-1	-2	-2	
W	-2	-2	-3	-2	-3	-4	-4	-3	-2	-3	-3	-3	-3	-2	-1	11	1	2	-2	-4	
F	-2	-3	-2	-2	-2	-3	-3	-3	-3	-3	-3	-1	0	0	0	1	6	3	-1	-4	
Y	-2	-3	-2	-2	-2	-2	-3	-2	-1	-2	-2	-1	-1	-1	-1	2	3	7	2	-3	
H	-3	-2	-2	-2	-1	1	-1	0	0	-1	0	-3	-3	-3	-2	-2	-1	2	8	-2	
P	-3	-2	-1	-1	-1	-2	-1	-1	-1	-1	-2	-2	-3	-3	-2	-4	-4	-3	-2	7	

# Search for sequence similarities: BLAST

Basic  
Local  
Alignment  
Search  
Tool



Algorithm for  
comparing a given  
sequence against  
sequences in a  
database



# BLAST algorithms

BLASTN

Your sequence  
QUERY

nucleotides

BLASTX

Translated DNA

BLASTP

Protein

TBLASTX

Translated DNA

TBLASTN

Protein



Databank  
SBJCT

nucleotides

Proteins

Proteins

Translated DNA

Translated DNA

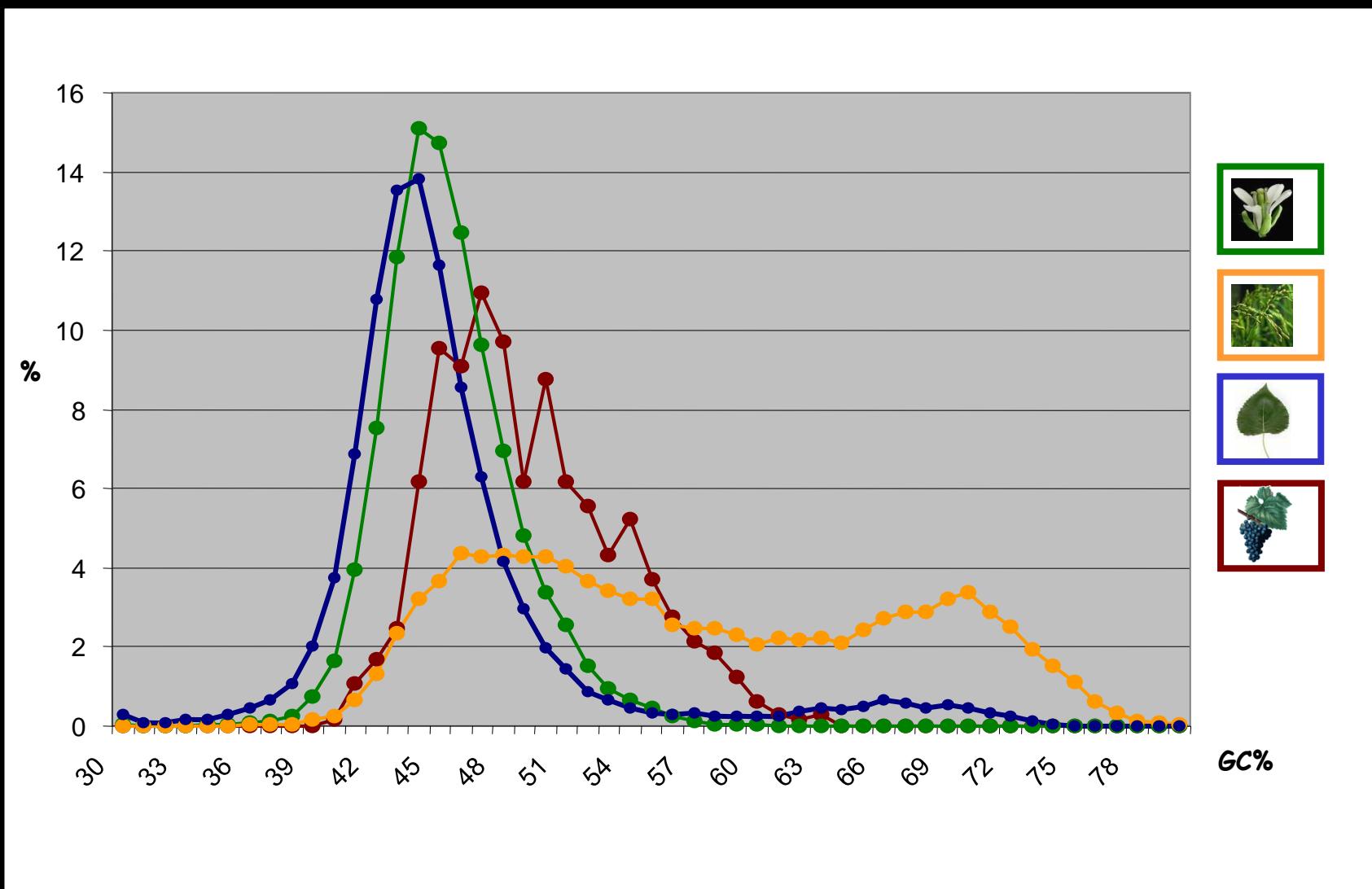
# The 'coding potential'

---

- Each organism use differently the genetic code redundancy
- Amino acids and synonym codons have not the same frequency in a genome (W rare, bias in codon usage)
- There are relations between successive codons.

The probability of the presence of a nucleotide depends to the previous K nucleotides (K: order of the Markov chain). This probability is different between coding and non-coding regions.

# Coding potential: learning step is necessary



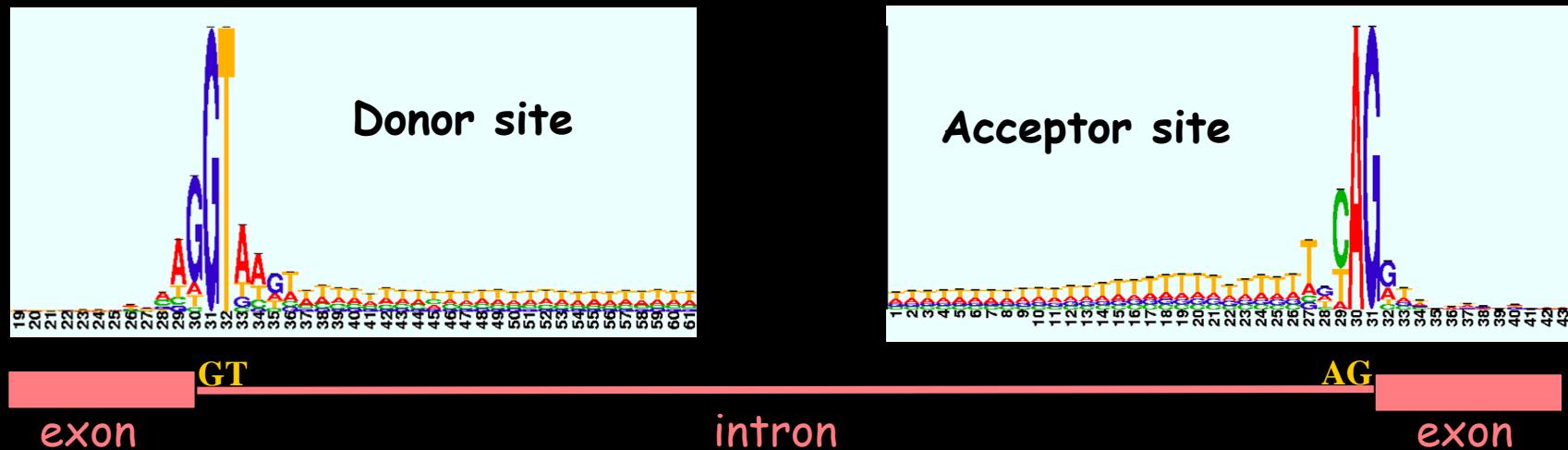
GC content in CDS

# Looking for splicing sites

Only in eukaryotic genomes...

## Positional Weight Matrices - WAM

LOGO



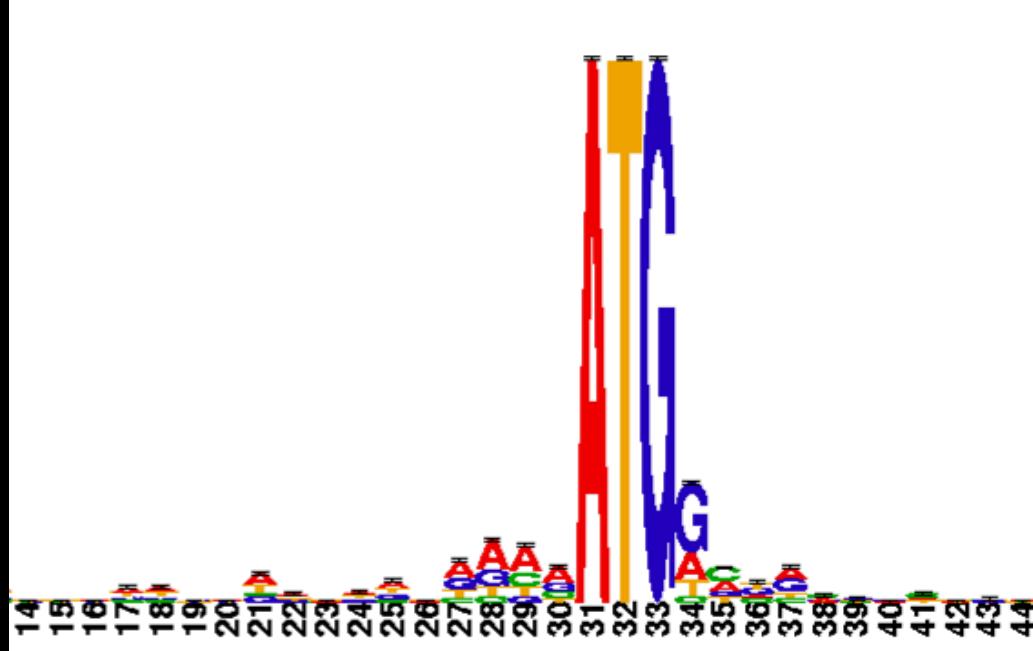
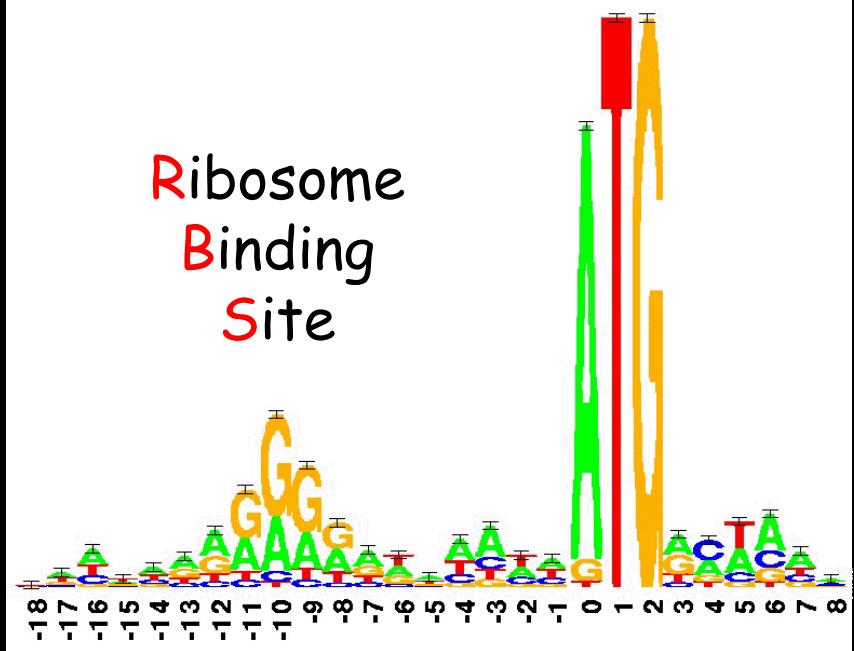
U2	[	GT::AG	99%
		GC::AG	1%
U12	[	AT::AT	
		AT::AC } 1%	

# Other signal: Translation initiation sites

Prokaryotes (*E. coli*)

Eukaryotes (*A. thaliana*)

Ribosome  
Binding  
Site

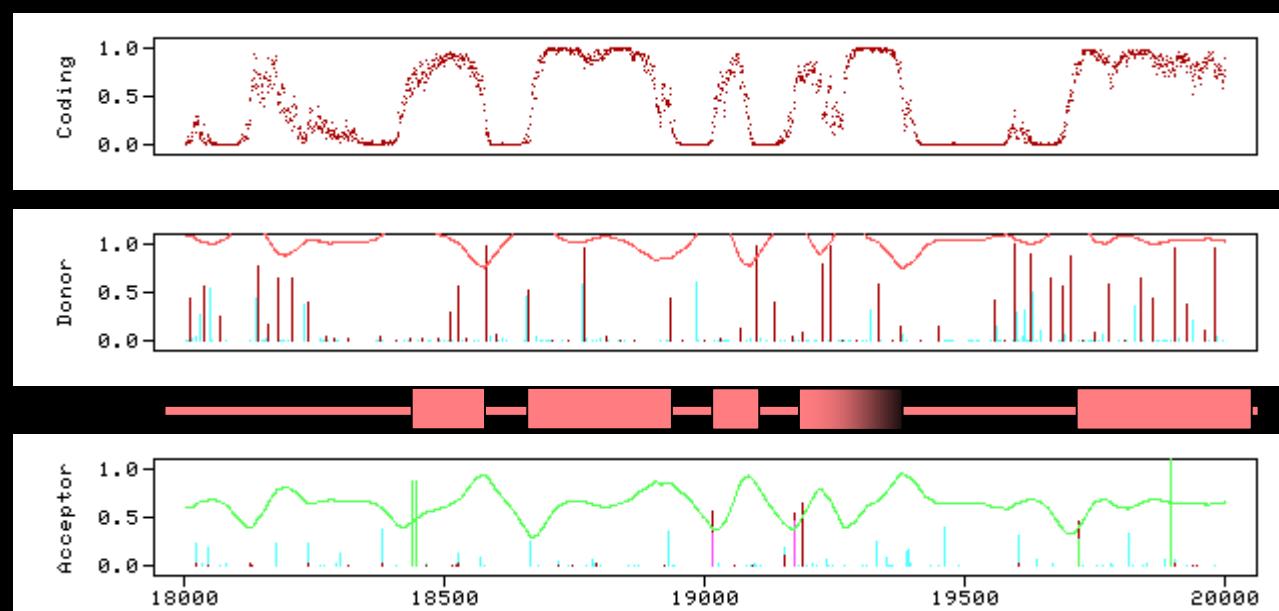


# Exon and intron predictions

## Restrictions

- minimum length of introns
- respect of reading frame
- initial exon: ATG codon
- terminal exon: STOP codon

NetGene2



The prediction only regards the coding exons  
The UTRs are difficult to predict

► CDS

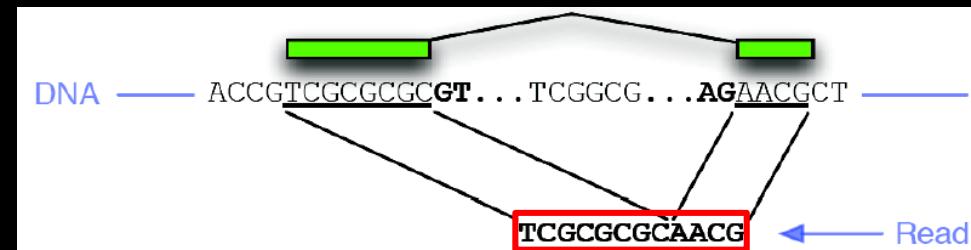
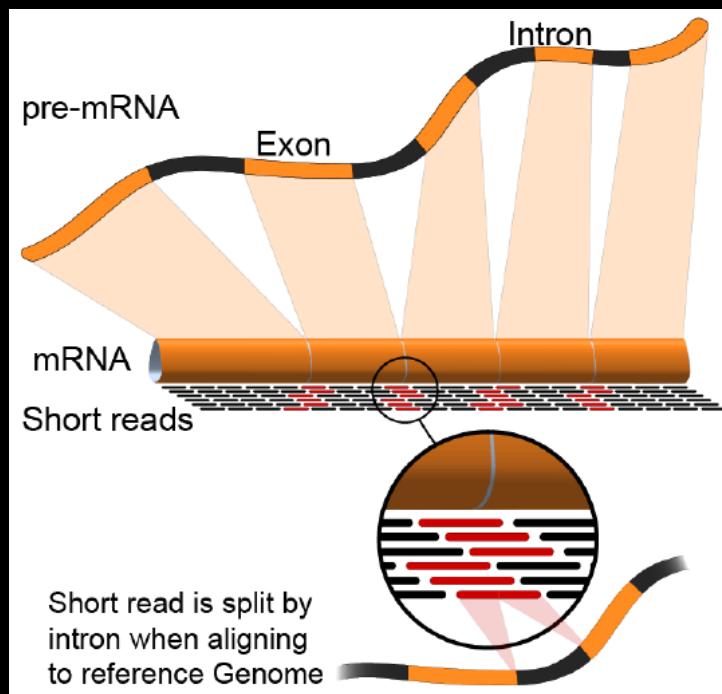
# Exon and intron predictions

ttcaactagca	acctcaaaca	gacaccatgg	tgcacactgac	tcctgaggag	aagtctgccg	Exon 1
ttactgcct	gtggggcaag	gtgaacgtgg	atgaagttgg	tggtgaggcc	ctggcagat	
tggtatcaag	gttacaagac	agguttaagg	agaccaatacg	aaactgggca	tgtggagaca	
gagaagactc	ttgggtttct	gataggcact	gactctctct	gcctatttgt	ctatttcccc	
acccttaggc	tgctgggtgg	ctacccttgg	accagaggt	tcttgagtc	ctttggggat	
ctgtccactc	ctgatgctgt	tatgggcaac	cctaagggtga	aggctcatgg	caagaaaatg	
ctcggtgcct	ttagtgtatgg	cctggctcac	ctggacaacc	tcaagggcac	ctttgccaca	Exon 2
ctgagtgagc	tgcactgtga	caagctgcac	gtggatctg	agaacttcag	ggtagtcta	
tgggaccctt	gatgtttct	ttccccttct	tttctatgg	taagttcatg	tcataggaag	
gggagaagta	acagggtaca	gtttagaatg	ggaaaacagac	gaatgattgc	atcagtgtgg	
aagtctcagg	atcgttttag	tttcttttat	ttgctgttca	taacaattgt	tttcttttgt	
ttaattcttg	ctttcttttt	ttttcttctc	cgcaattttt	actattatac	ttaatgcctt	
aacattgtgt	ataacaaaag	gaaatatctc	tgagatacat	taagtaactt	aaaaaaaaac	Exon 3
tttacacagt	ctgcctagta	cattactatt	tggaatatat	gtgtgcttat	ttgcatattc	
ataatctccc	tactttattt	tcttttattt	ttaattgata	cataatcatt	atacatattt	
atgggttaaa	gtgtaatgtt	ttaatatgtg	tacacatatt	gaccaaatac	gggtaatttt	
gcatttgtaa	ttttaaaaaa	tgctttcttc	ttttaatata	ctttttgtt	tatcttattt	
ctaatacttt	ccctaattctc	tttctttcag	ggcaataatg	atacaatgta	tcatgcctct	
ttgcaccatt	ctaaagaata	acagtataaa	tttctgggtt	aaggcaatag	caatatttct	Exon 3
gcatataaat	atttctgcat	ataaaattgt	actgatgtaa	gaggtttcat	attgctaata	
gcagctacaa	tccagctacc	attctgcttt	tattttatgg	ttgggataag	gctggattat	
tctgagtcca	agctaggccc	tttgctaat	catgttcata	ccttttatct	tcctcccaca	
gctcctgggc	aacgtgctgg	tctgtgtgct	ggcccatcac	tttggcaaag	aattcaccgg	
accagtgcag	gctgcctatc	agaaaagtgg	ggctgggtgt	gctaattgccc	tggcccacaa	
gtatcactaa	gctcgctttc	ttgctgtcca	atttcttatta	aaggttcctt	tgttccctaa	Exon 3
gtccaactac	taaactgggg	gatattatga	agggccttga	gcatctggat	tctgcctaata	
aaaaaaacatt	tattttcatt	gcaatgtatgt	attnaaatta	tttctgaata	ttttactaaa	

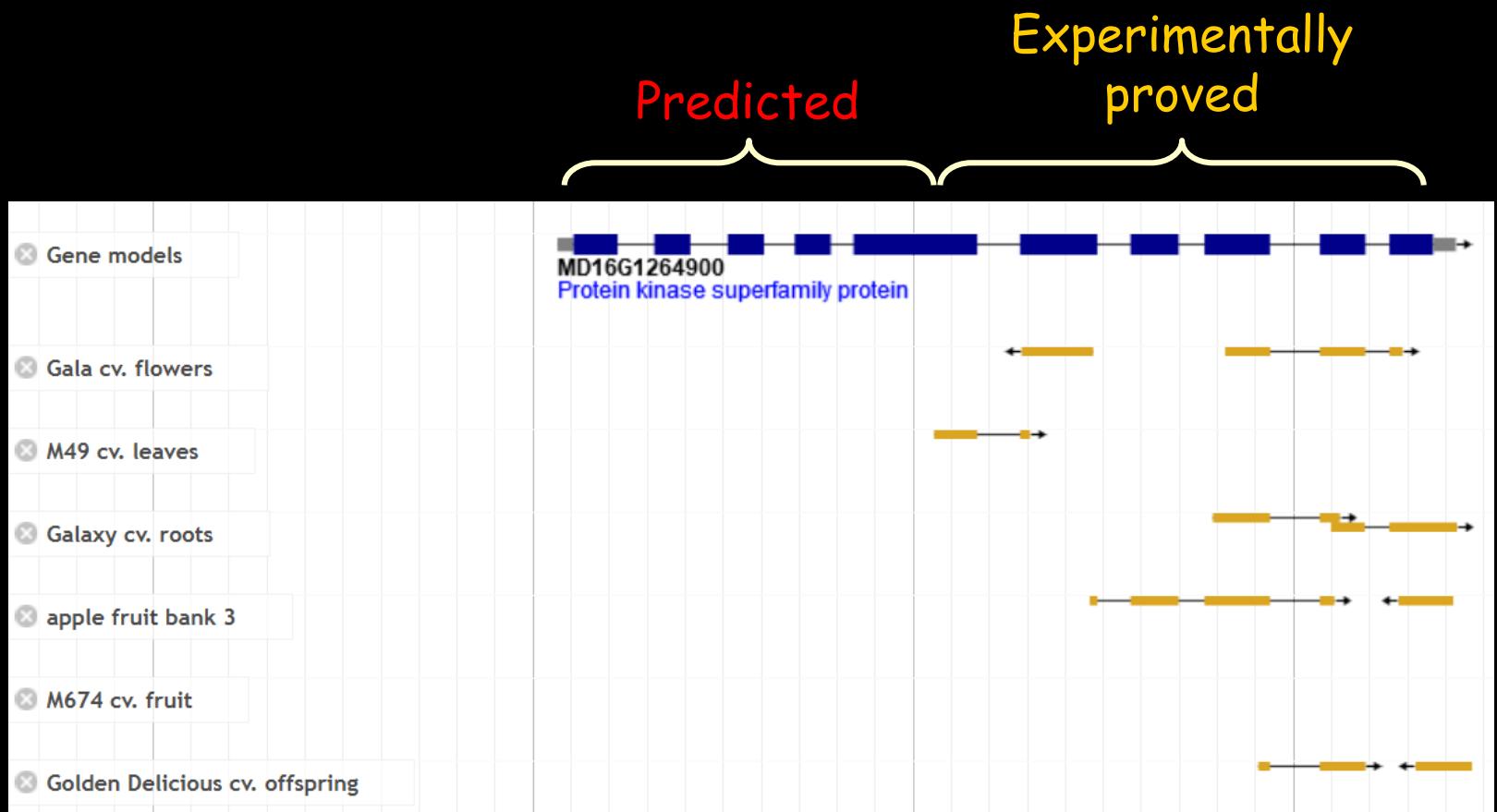
# Transcripts and genomic sequences

Alignment of cDNA, EST or consensus of short read contigs with the cognate genomic sequence (gmap, star, sim4, spidey...) :

- Spliced full length alignment
- Deal with low quality sequences
- Take into account known splicing sites

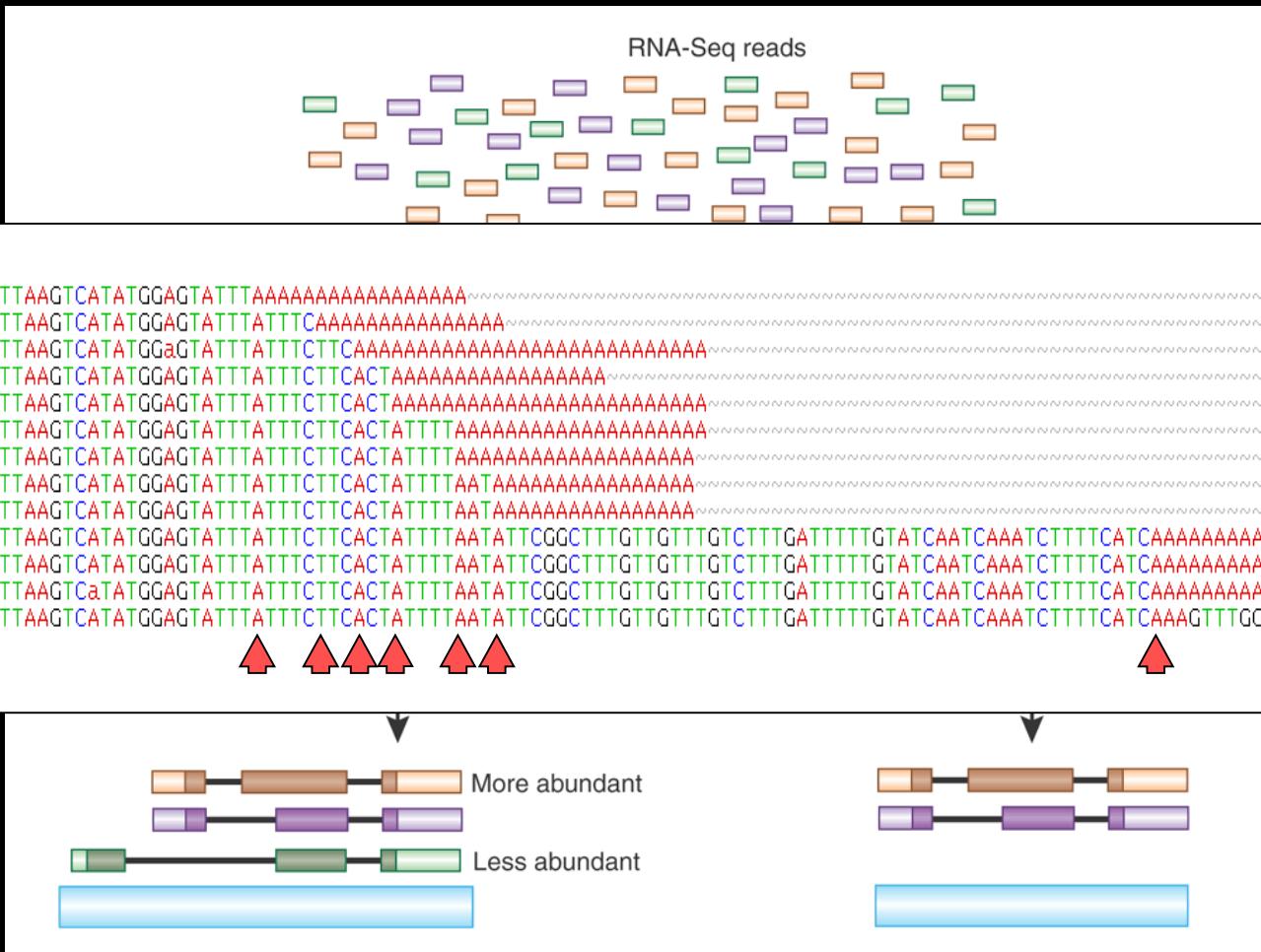


# Transcripts and genomic sequences



Heterogeneous quality of structural annotation inside a gene !

# Transcript mapping / de novo assembly



# Assembly goals

- more informative consensus sequence
  - reduction of complexity
  - remove sequencing errors
  - highlight gene structure and alternative events

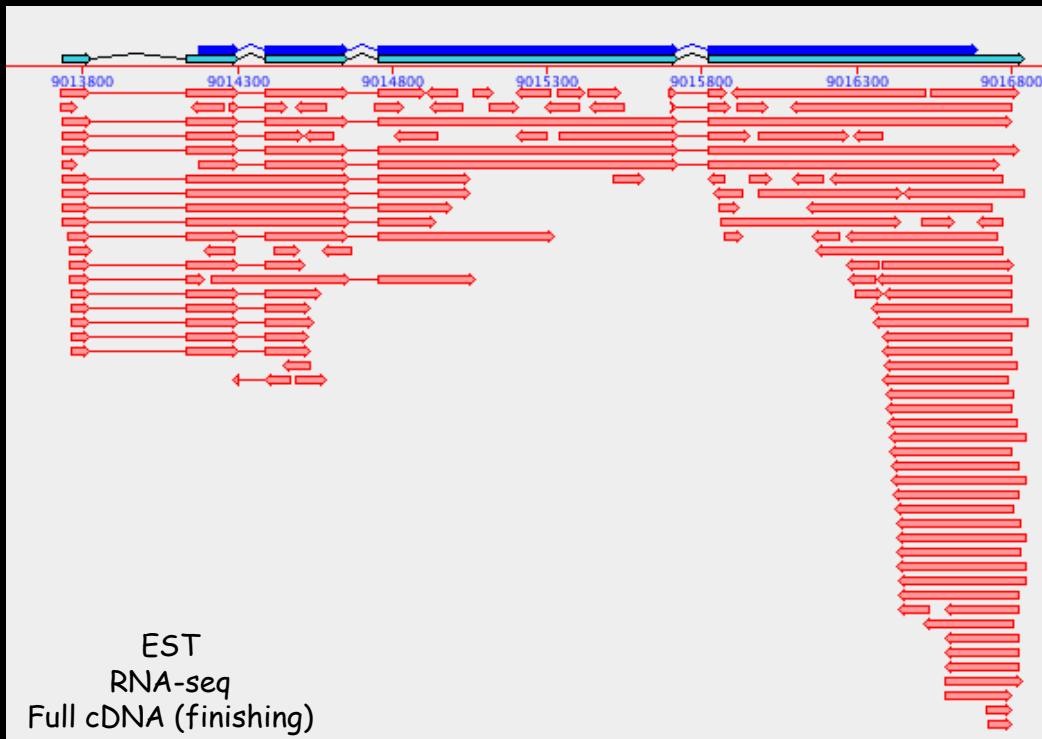
# Transcript exploitation

## Qualitative

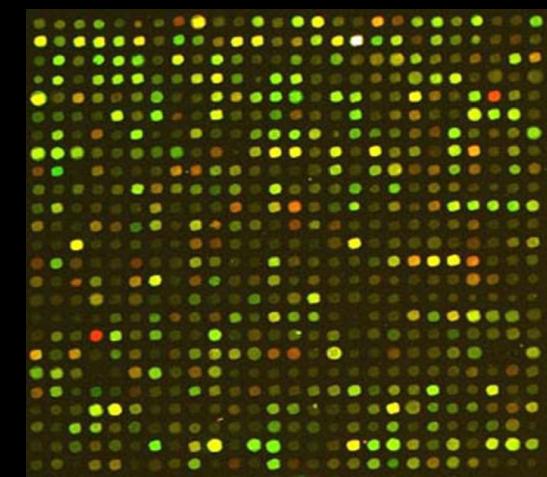
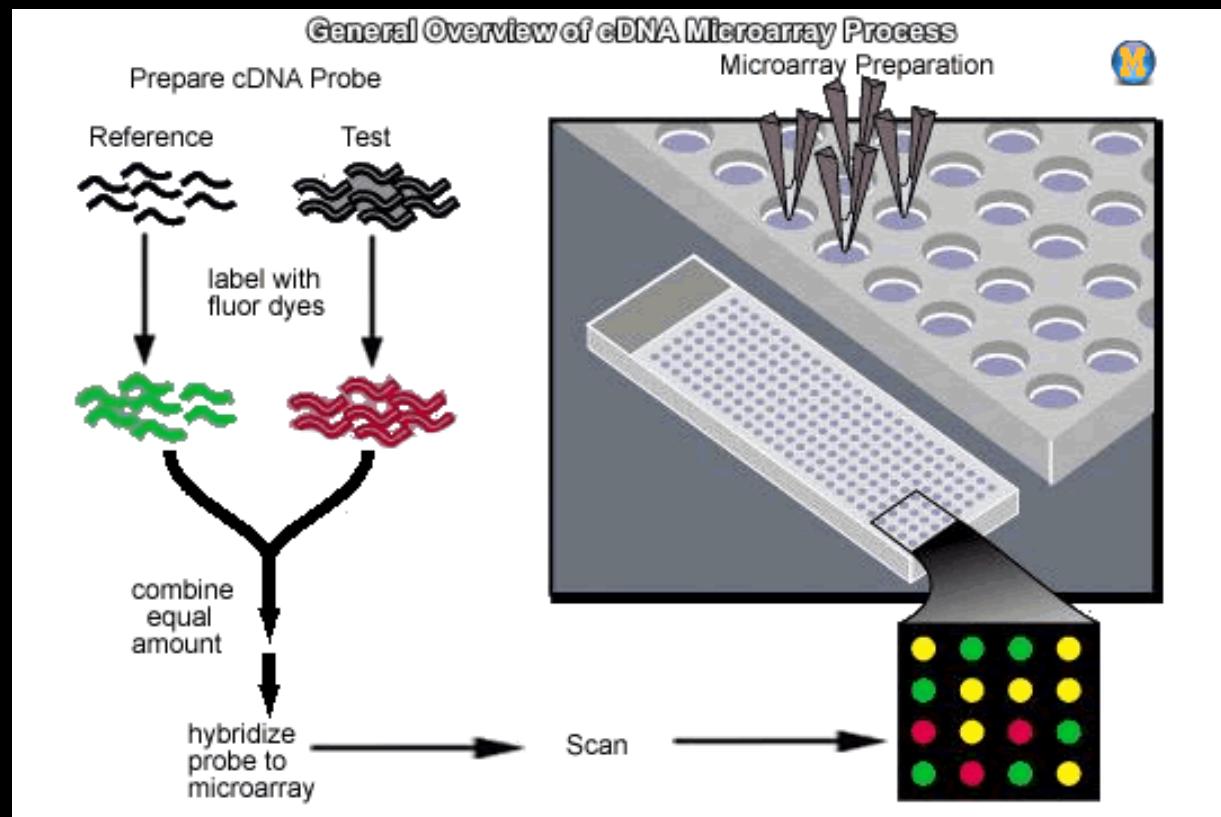
Proof of gene transcription in a specific condition / intron-exon structure / UTRs

## Quantitative

The number of cognate reads reflects the relative proportion of transcripts in the biological sample used



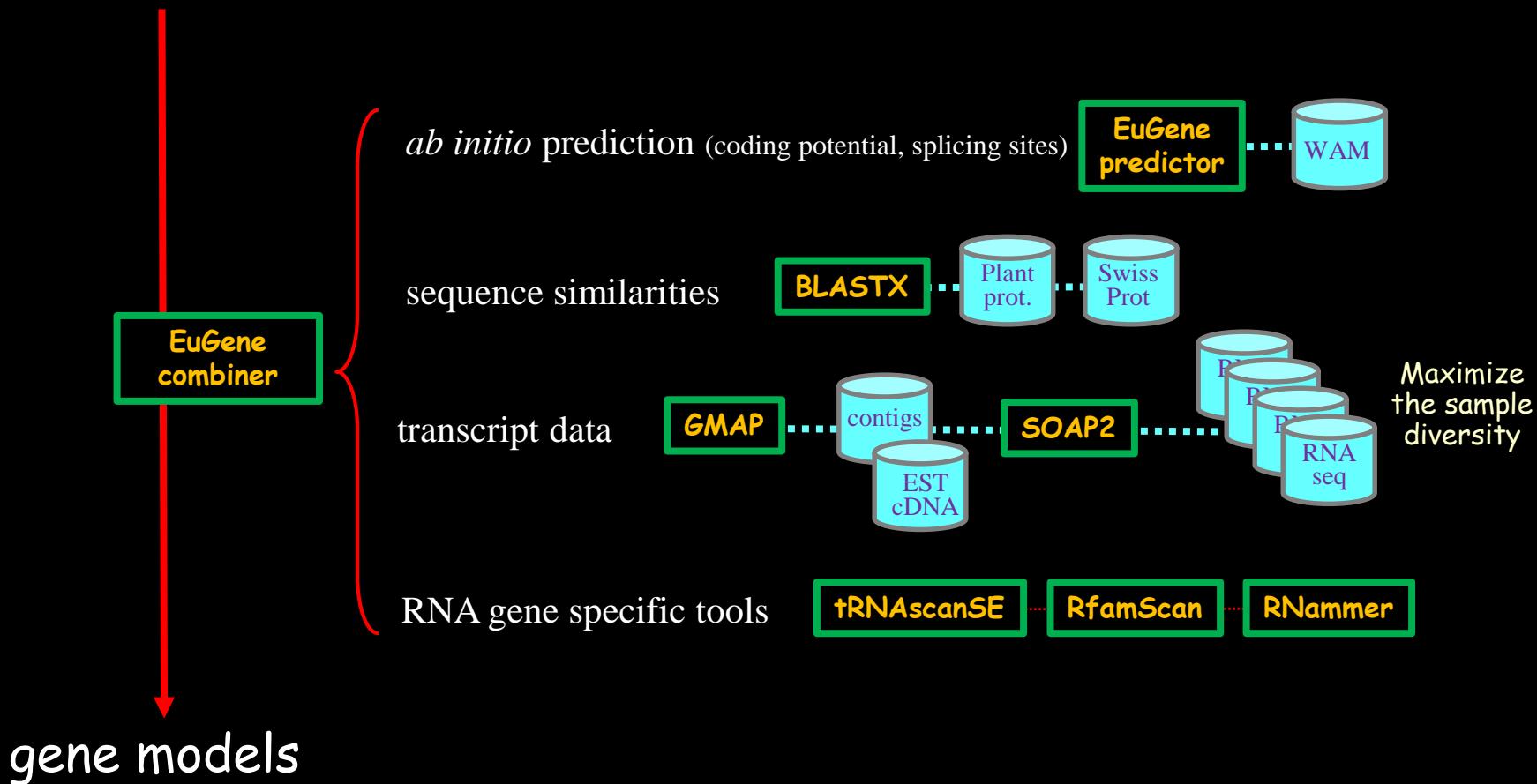
# Expression analysis with micro-array



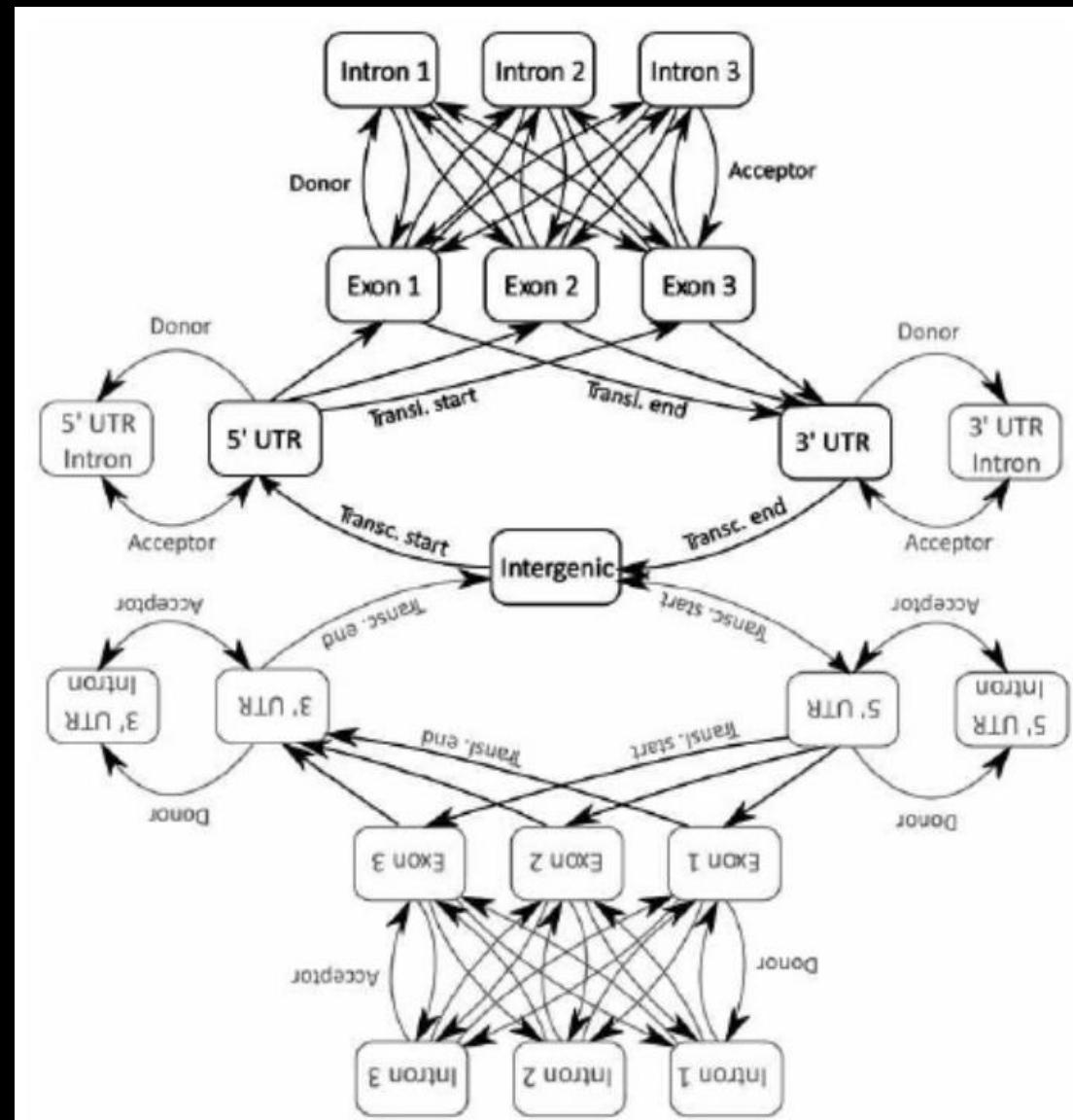
# EuGene: a complete pipeline for gene modeling

```
ATAACATGATCATGTTAACCGTTGAATC  
AAAATGATTAAAAGAACATAATGTCGTG  
ATGTCGTTAAATCTCAGAAGTGTGG  
AGAGAGGATTCAAGGGTAACCATAACCCA  
TGAGGCTGGTAAATGGTACTATCAGCT  
ACATTGACGCAGGATCCAAGCAAGCTTAT  
TCAGCTTGACCTGCATCACCGTTATTCA  
TGACAACCTGGAGGAAGCCGATACAGACTA  
CTGAGAAGTTGGAATGCTTAGACTGGT
```

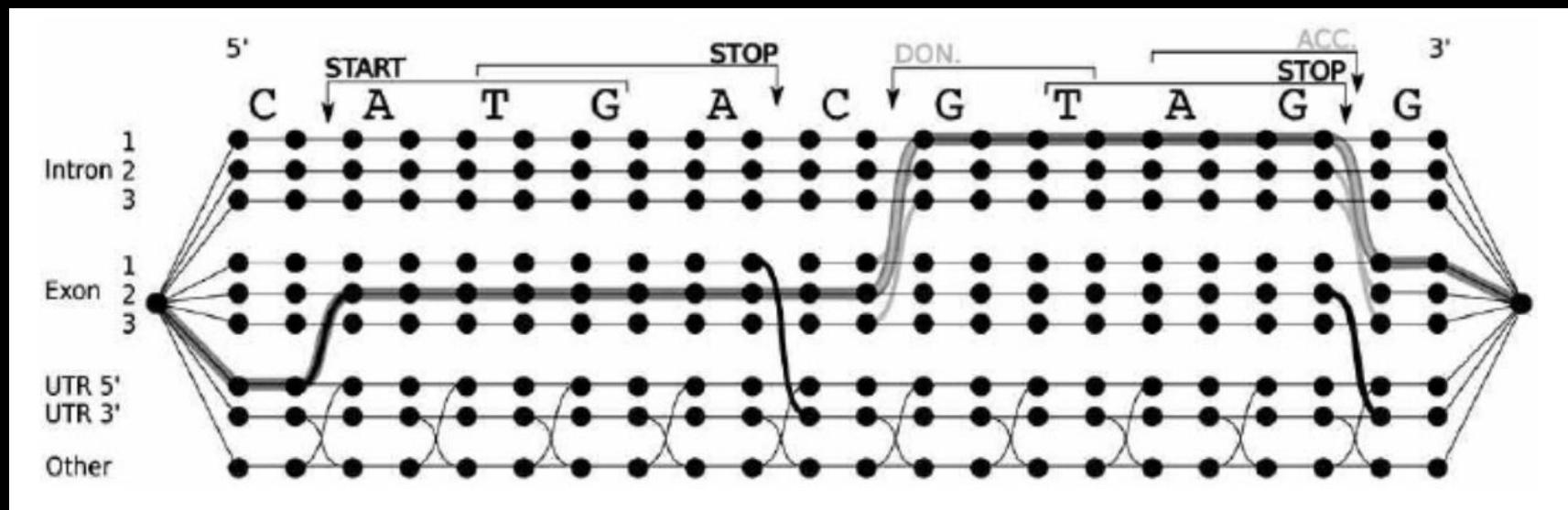
chromosomes, contigs



# Gene modeling : the states of a genomic sequence

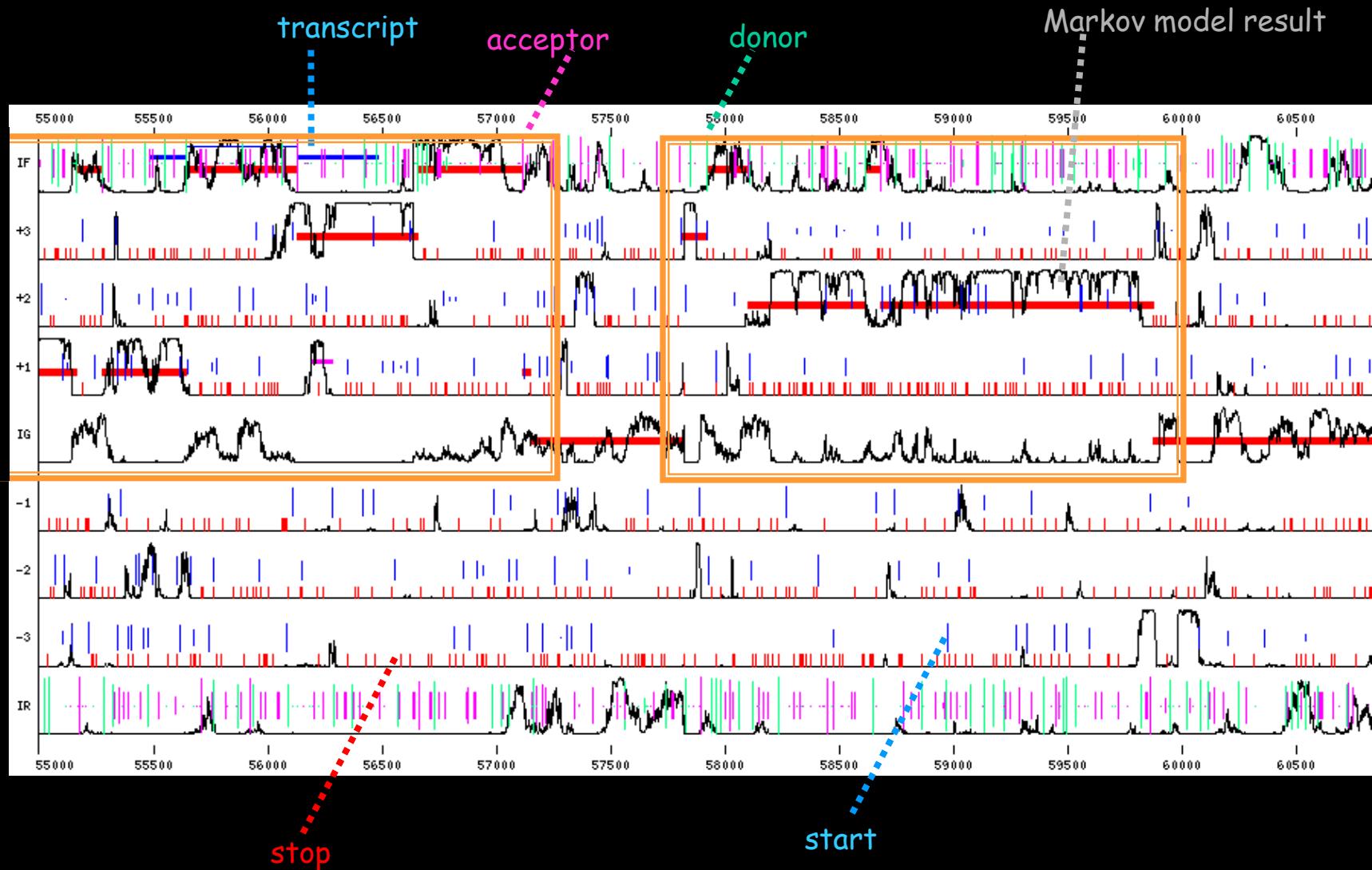


# Gene modeling : the states of a genomic sequence



Eugene (T. Schieex)

# Integration for gene prediction



# Examples of gene prediction tools

---

GeneMark.HMM

Glimmer

GeneFinder

EuGene

FgeneSH

Augustus

Genie

HMMgene

MagPie

GeneId

Grail

GeneScan

GeneWise

...

~ 70 predictors

# Problems and errors: where and why ?

---

## Gene extremities are difficult to predict

- ✓ No or low differences between intergenic and UTR regions
- ✓ in 3': polyadenylation signals are rare or degenerated (AAUAAA)
- ✓ in 5': no systematic CpG islands and/or TATA box
- ✓ intergenic regions can be short and introns can be very long
- ✓ N- and C-terminal regions are often lowly conserved

## ► Gene breaking and gene merging

## Bottlenecks of spliced alignments

- ✓ Assembly step can produce chimeric transcripts
- ✓ Very short exons can fail the alignment process (<10 bp)
- ✓ Mainly GT::AG sites are looked for to open gaps
- ✓ Sequencing of partially spliced transcripts

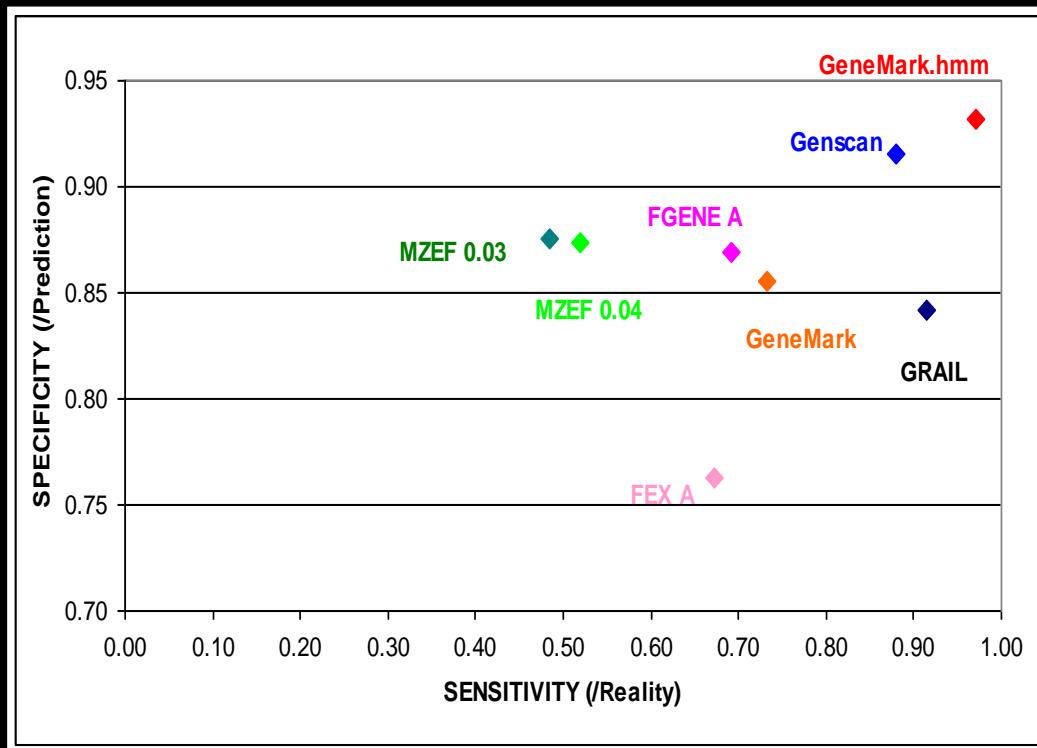
# Evaluation and efficiency of gene predictors

true positive (TP): reality well predicted

true negative (TN): false predict to be false

false positive (FP): false predict to be true

false negative (FN): reality not predicted



Sensitivity

$$\frac{TP}{TP+FN}$$

Fraction of reality  
that is  
well predicted

Specificity

$$\frac{TP}{TP+FP}$$

Fraction of the  
prediction that is  
true

...at different levels: splicing sites, exons or genes

# Evaluation of genome annotation



Benchmarking Universal Single-Copy Orthologs  
Plant dataset : 1440 genes

Sunflower : 92%

Medicago : 95,4%

Rosa : 96,5%

Apple : 96,8%



complete	<b>BUSCO</b> v3	96,8 %
fragmented		1,3 %
missing		1,9 %
transcript support		89,8 %
PFAM signature		74,8 %
TAIR or SwissProt homolog		89,6 %

} U : 93,2 %

▼ expertise of 1194 genes (50 families)

CDS ok	site	ATG	exons	split	merge	underpred	pseudo
<b>753</b>	<b>18</b>	<b>55</b>	<b>54</b>	<b>23</b>	<b>20</b>	<b>32</b>	<b>239</b>
<b>78,8 %</b>	<b>1,9 %</b>	<b>5,8 %</b>	<b>5,7 %</b>	<b>2,4 %</b>	<b>2,1 %</b>	<b>3,4 %</b>	<b>20 %</b>

955 functional genes

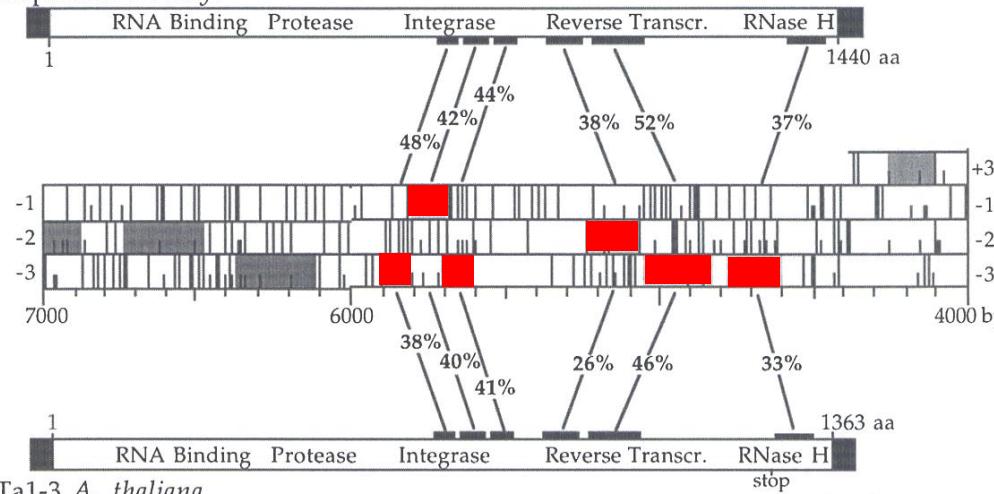
# Annotation of transposable elements

## RepeatMasker

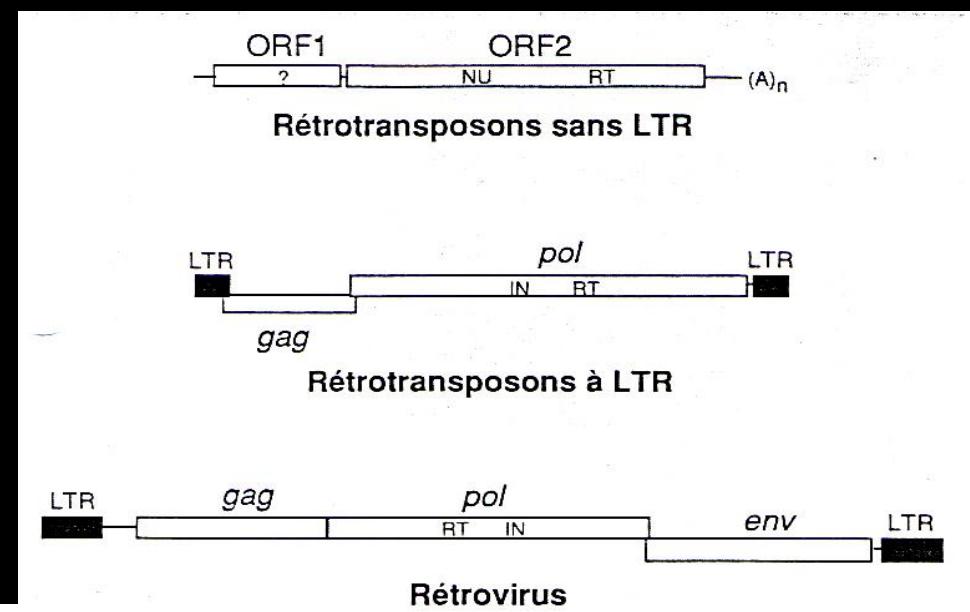
- repeated sequences
- Transposables elements  
(RepBase)

## Repet pipeline

Hopscotch *Z. mays*

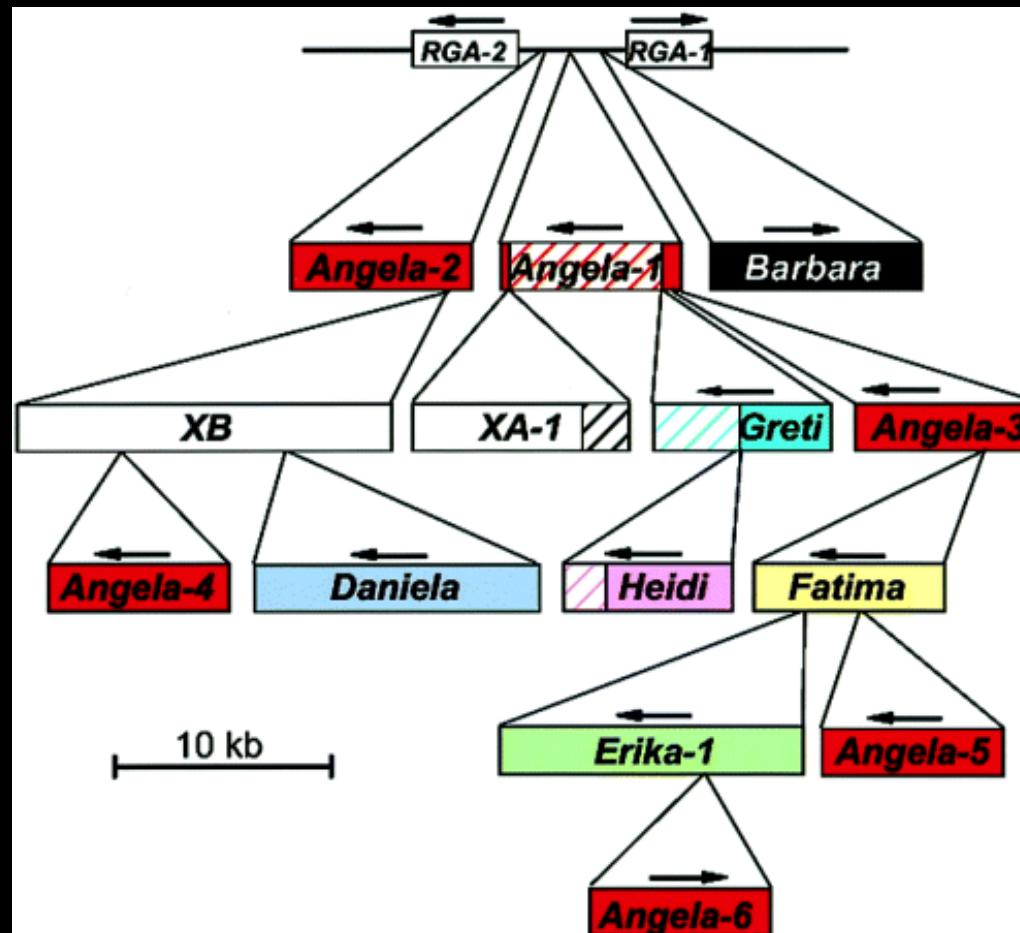


Ta1-3 *A. thaliana*



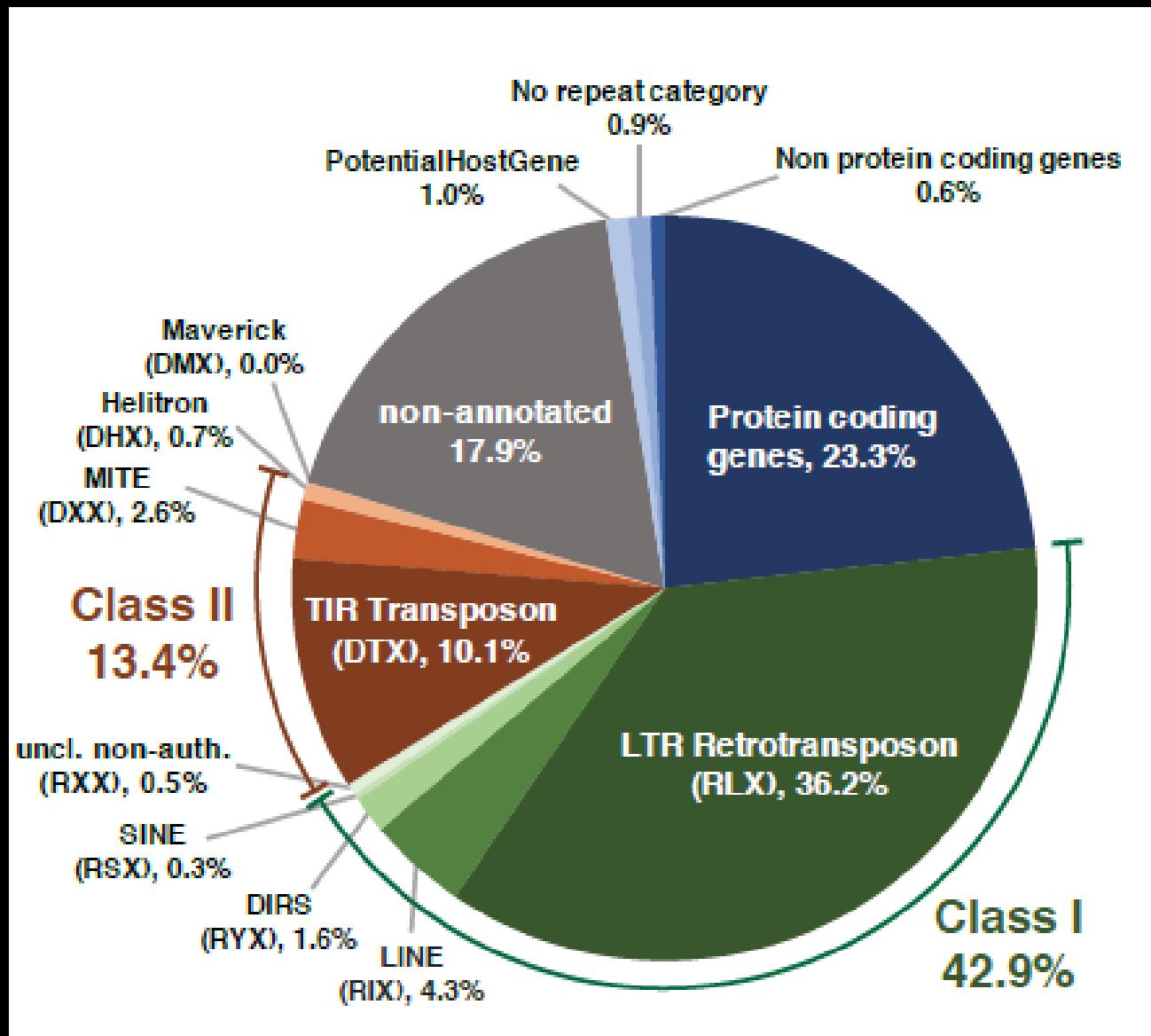
Quesneville *et al.*  
Combined evidence annotation of transposable elements in genome sequences.  
*PLoS Comput Biol.* 2005 Jul;1(2):e22  
PMID: 16110336

# Annotation of transposable elements: evolution



Wicker et al. (wheat, 2001)

# TE overview in the apple genome



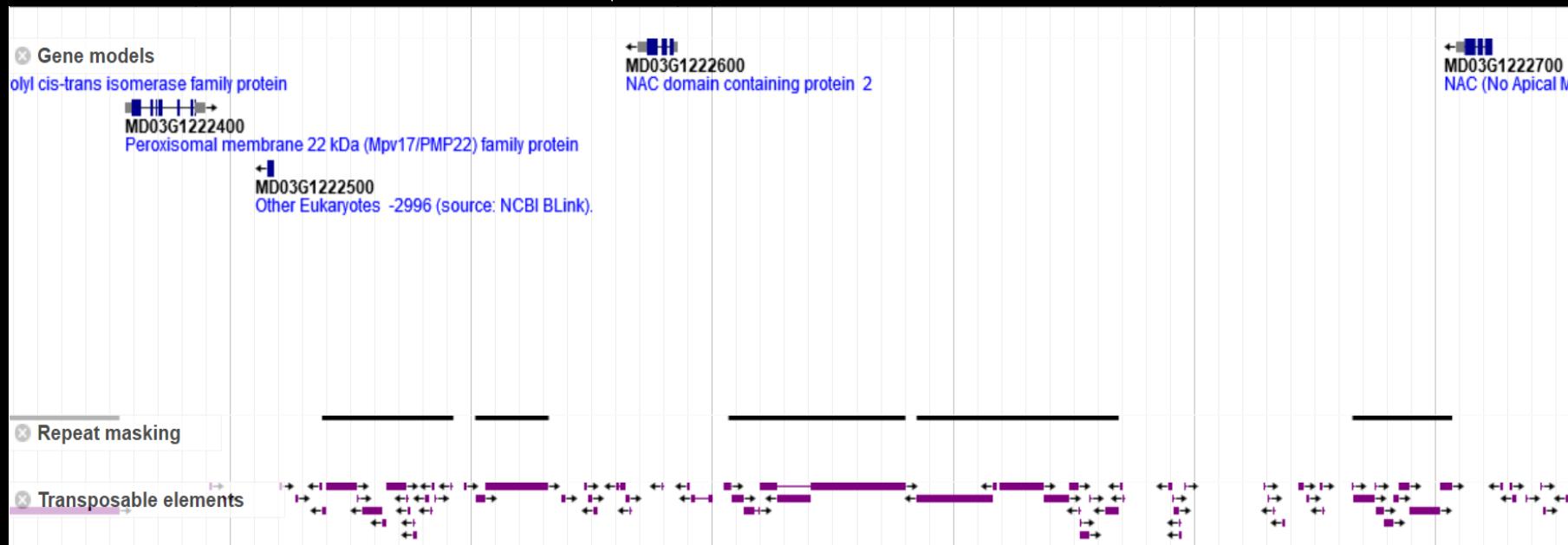
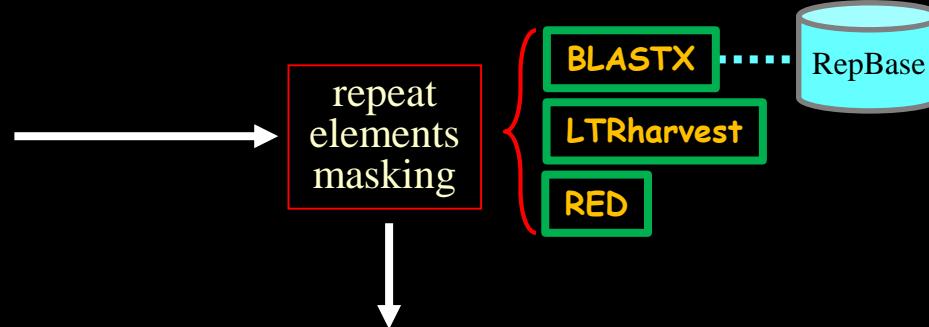
# Repeat masking

Goals : to avoid over-prediction of genes in

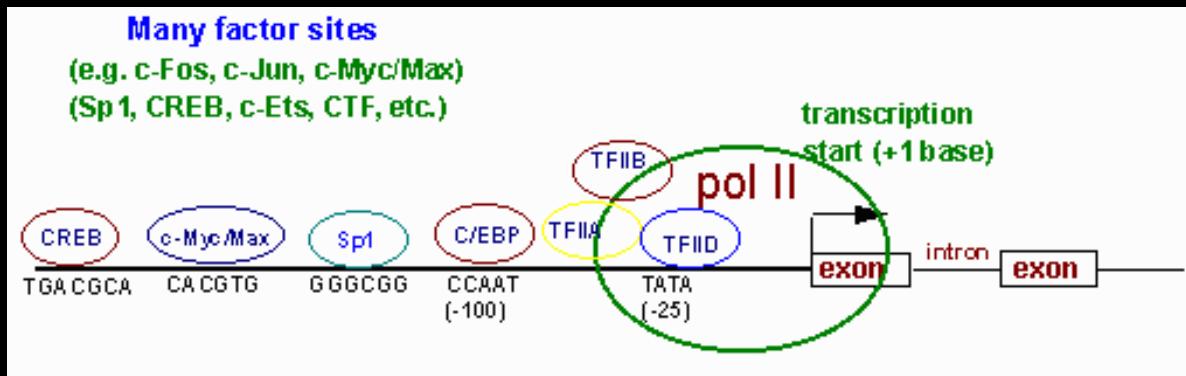
- disrupted coding regions of TE (complete or not)
- low complexity regions

genome

```
ATAACATGATCATGTTAACGTTGTAATC  
AAAATGATTAAAAAGAACTAATGTCGG  
ATGTTGCTTAAATCTCAGAAGTGG  
AGAGAGGATTCAGGGTAAACCATAACCA  
TGAGGCTGTGAAATGGTACTATCAGCT  
ACATTGACCCAGGATCCAAGCAAGCTTAT  
TCAGCTTGACCTGCATCACCGTTATTOA  
TGACAACCTGGAGAACCGGATACAGACTA  
CCTGAGAAGTTGGAATGTCTTAGACTGGT
```

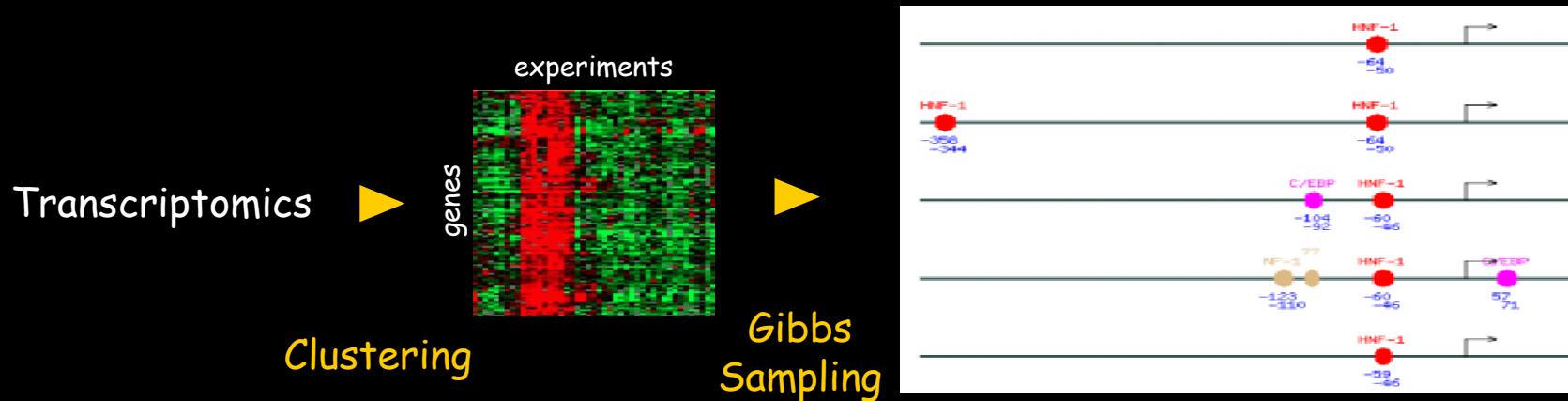


# Annotation of promoters



- ✓ short and degenerated motifs ► numerous false positives
  - ✓ Transcription Start Sites (TSS) identification are based on transcripts
  - ✓ No systematic CpG islands and/or TATA box
- Never done at the genome scale

# Annotation of promoters in co-expression clusters



- TRANSFAC
- MOTIF SAMPLER (TOUCAN package)
- RSAT
- MATINSPECTOR
- VISTA
- PROMO
- ALIGNACE
- MEME
- SCPD (yeast)
- PLACE (plant)
- PLANTCARE (plant)

# Example of plant genomes



## *ab initio*

Genscan+  
GeneMark.hmm  
GlimmerA  
GeneSplicer

Genscan+  
GeneMark.hmm  
GlimmerM  
FGENESH

EUGENE  
GrailExp6  
GeneWise  
FGENESH

GeneId  
GlimmerHMM  
SNAP  
EUGENE

FGENESH  
EUGENE

EUGENE

## similarities

BLAST

BLAST

BLAST  
S&W

EXOFISH  
BLAST

BLAST

BLAST

## transcripts

PASA  
BLAT  
SIM4

PASA2

SIM4

EST2Genome  
GeneWise

PASA

SOAP2  
GMAP

## combiner

AAT  
COMBINER

AAT  
COMBINER

JGI R.Syst.

GAZE

EUGENE

EUGENE

TIGR  
TAIR

TIGR

JGI  
PSB  
ORNL

GENOSCOPE  
URGV

IMGAG

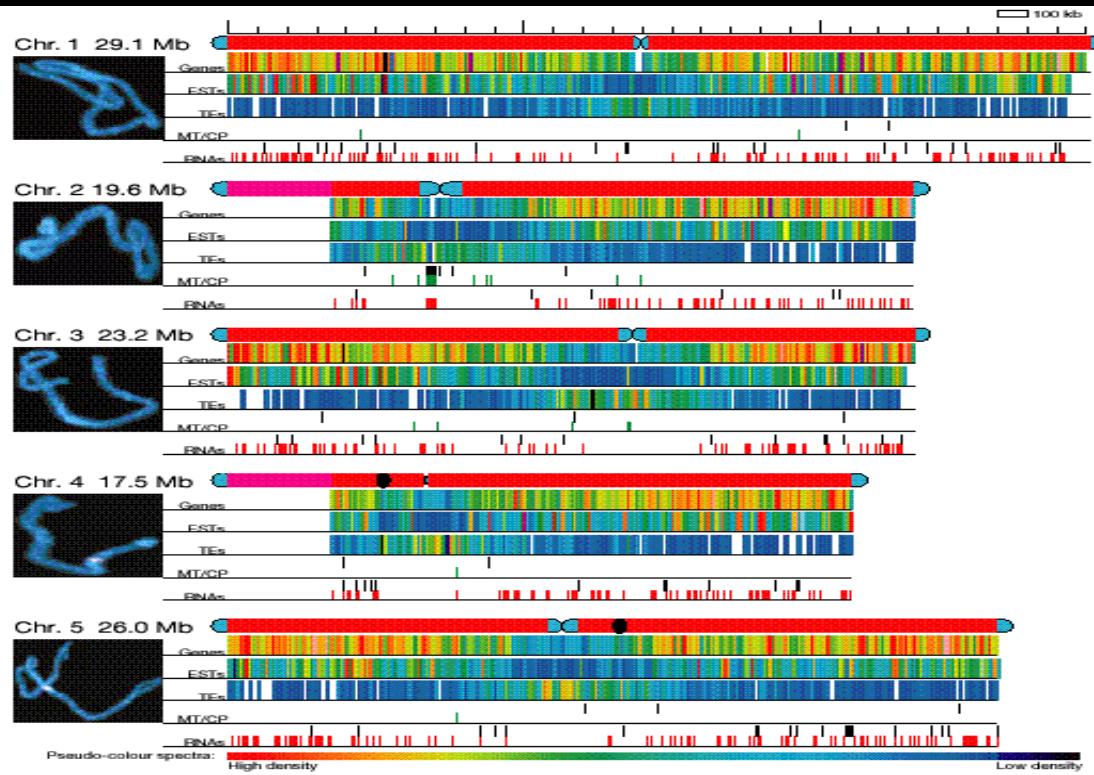
IRHS

# Conclusions about structural annotation

---

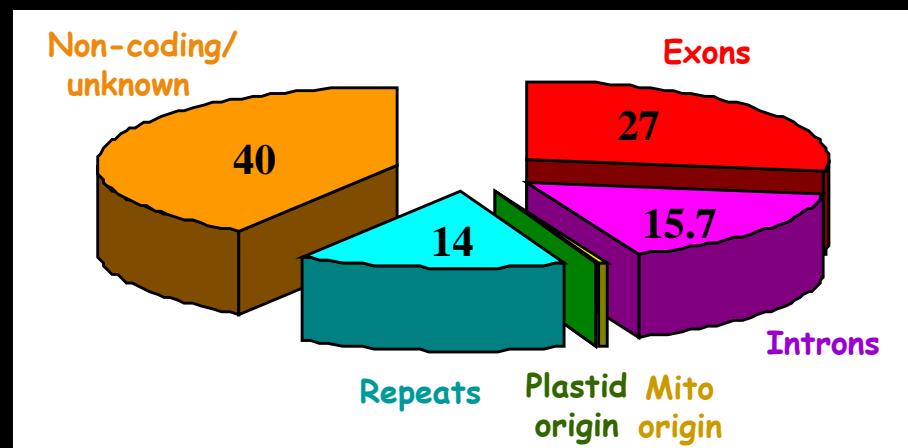
- Automatic annotation is quick and provides (preliminary) data wanted by the scientific community
- Numerous errors but not too much missed full coding genes
- The different prediction pipelines have to be known, understood and took into account to read and exploit the annotations
- Other approaches are necessary (by gene family, by experts..)
- Take precautions:
  - BLASTX, P or PSI (bank of proteins + translations of CDS)
  - Prodom, Pfam...
- Mix of proved and only supposed and predicted data
- Rare and novel situations can not be predicted !

# The Arabidopsis genome

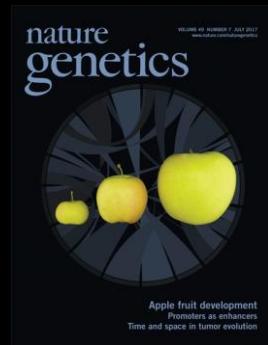


118 997 677 bp  
27 665 genes  
4 853 pseudogenes + TE

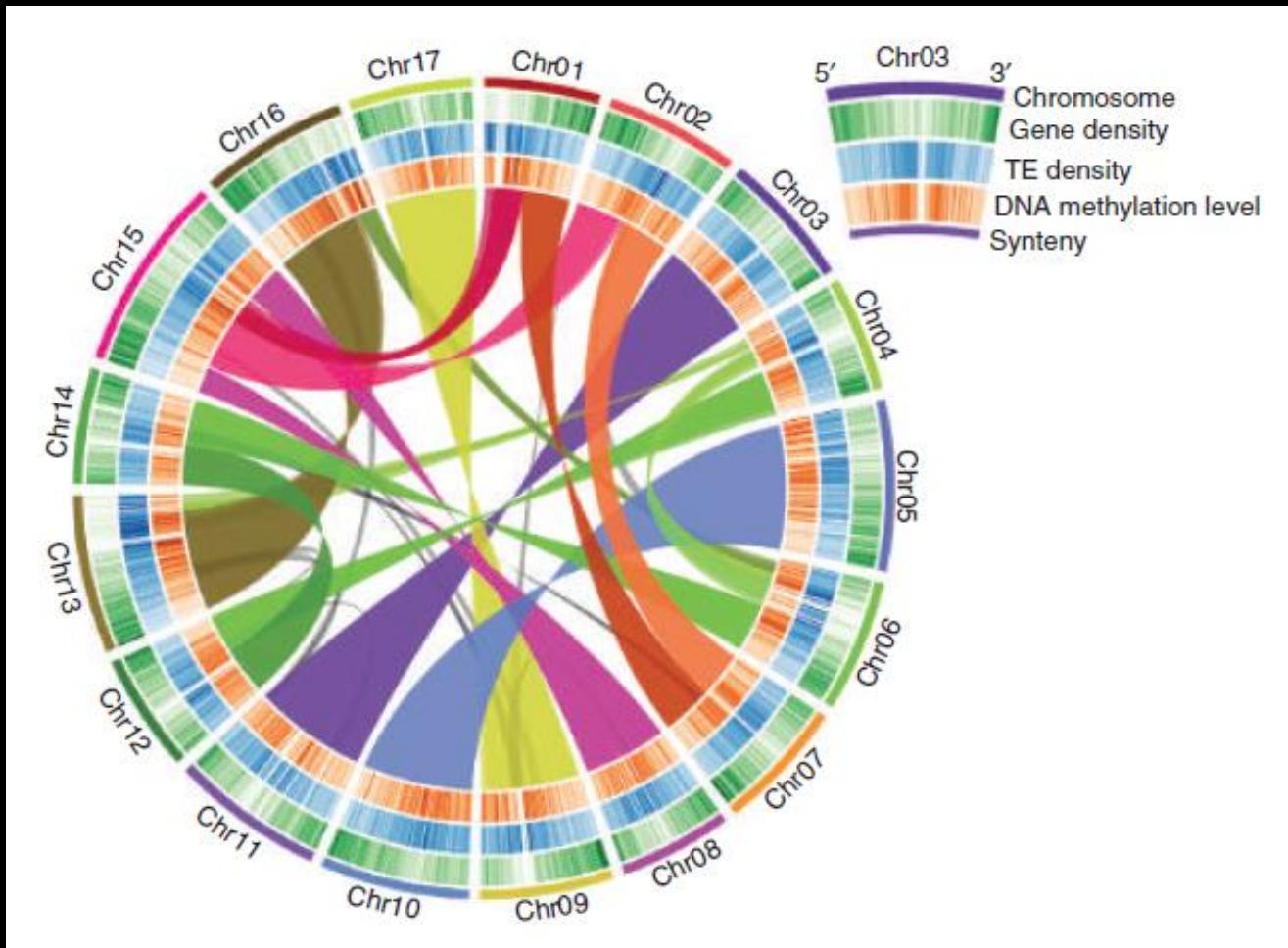
10 695 alternatively spliced genes



# The *Malus domestica* genome



<b>coding genes</b>	<b>45116</b>
snoRNA	410
rRNA	4369
tRNA	654
miRNA	141
snRNA	181
other ncRNA	54
undefined ncRNA	1816
<b>total genes</b>	<b>52741</b>



# Format of annotations @ GB/EMBL

---

```
repeat_region 11095..11216
  /rpt_family="(TAAAA)n"
gene    12166..13245
  /gene="T27C4.4"
  /note="similar to ribosomal protein L17 GB:AAA34113.1 from
[Nicotiana tabacum]"
mRNA   join(12166..12188,12416..12700,12785..12826,13063..13245)
  /gene="T27C4.4"
CDS     join(12176..12188,12416..12700,12785..12826,13063..13145)
  /gene="T27C4.4"
  /codon_start=1
  /product="ribosomal protein L17, putative"
  /protein_id="AAF63771.1"
  /db_xref="GI:7547099"
  /translation="MSKRGRRGGTSGNKFRMSLGLPVAATVNCADNTGAKNLYIISVKG
IKGRLNRLPSACVGDMVMATVKKGPDLRKVLPAVIVRQRPWRRKDGVFMYFEDNA
GVIVNPKGEMKGSAITGPIGKECADLWPRIASAANAIV"
gene    <13674..>14042
  /gene="T27C4.5"
  /note="predicted by genscan+"
mRNA   join(<13674..13775,13914..>14042)
  /gene="T27C4.5"
CDS     join(13674..13775,13914..14042)
  /gene="T27C4.5"
  /note="hypothetical protein"
  /codon_start=1
  /protein_id="AAF63772.1"
  /db_xref="GI:7547100"
  /translation="MVNPVGFRFRPTKEEIVDHYLRPTNFDGDTSHVDRNIMFMQDNR
NDYRPPNSLTGVFSDCSSDDNDSLLSPKTVS"
repeat_region complement(15431..15636)
  /note="Limpet1 transposon GB U76697"
```

# Format of annotations @ GB/EMBL

repeat\_region 11095..11216

/rpt\_family="(TAAAAA)n"

gene 12166..13245

/gene="T27C4.4"

/note="similar to ribosomal protein L17 GB:AAA34113.1 from  
[Nicotiana tabacum]"

mRNA join(12166..12188,12416..12700,12785..12826,13063..13245)

/gene="T27C4.4"

CDS join(12176..12188,12416..12700,12785..12826,13063..13145)

/gene="T27C4.4"

/codon\_start=1

/product="ribosomal protein L17, putative"

/protein\_id="AAF63771.1"

/db\_xref="GI:7547099"

/translation="MSKRGRRGGTSGNKFRMSLGLPVAATVNCADNTGAKNLYIISVKG"

IKGRLNRLPSACVGDMVMATVKKGPDLRKVLPAVIRQRKPWRRKDGVFMYFEDNA

GVIVNPKGEMKGSITGPIGKECADLWPRIASAANAIV"

gene <13674..>14042

/gene="T27C4.5"

/note="predicted by genscan+"

mRNA join(<13674..13775,13914..>14042)

/gene="T27C4.5"

CDS join(13674..13775,13914..14042)

/gene="T27C4.5"

/note="hypothetical protein"

/codon\_start=1

/protein\_id="AAF63772.1"

/db\_xref="GI:7547100"

/translation="MVNPVGFRFRPTKEEIVDHYLRTNFQDTSHVDRNIMFMQDNR"

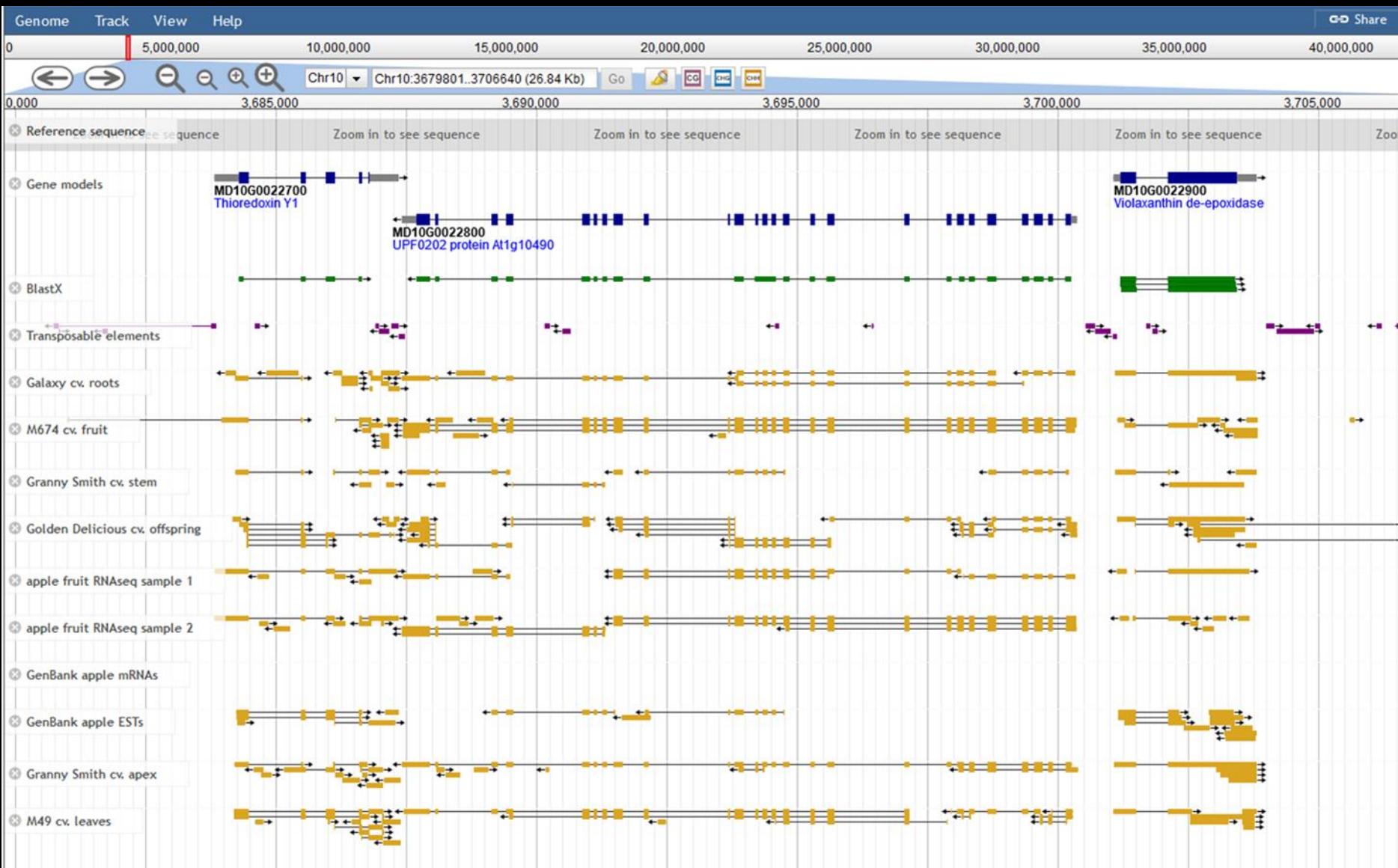
NDYRPPNSLTGVFSDCSSDDNDSLLSPKTTS"

repeat\_region complement(15431..15636)

/note="Limpet1 transposon GB U76697"

Feature  
/Qualifier

# Genome Browser



# Annotation tool

File Select View Goto Edit Create Write Run Display

Selected feature: bases 6 (/colour=6 /label=\* /note="gtattt, splice donor sequence")

EMBL Entry: c29A3.seq c29A3.tab

24000 24800 25600 26400 27200 28000 28800 29600 30400 31200

ur SPBC29A3.15c trt1 SPBC29A3.13c

G K V V F \* V S I L T L F \* K Y + L F A K N T L K Y L V Y F F K R  
A K + F S E Y P S # H C S E N I S F L Q K I L # N T S F I F S N  
Q S S F L S I H P N I V L K I L A F C K K Y F K I P R L F F Q T  
GCAAAGTAGTTTCTGAGTATCCATCCTAACATTGTTCTGAAAAATATTAGCTTTGCAAAAAACTTTAAACCTCGTTATTTTCAACG  
27520 27530 27540 27550 27560 27570 27580 27590 27600 27610  
CGTTTCATCAAAAGACTCATAGGTAGGATTGTAACAAGACTTTATAATCGAAAAACGTTTTATGAAAATTTATGAAAGTTATGCAACG  
L T T K Q T D M R V N N Q F Y # S K A F F V K F Y R T # K K L R  
C L L K R L I W G L M T R F I N A K Q L F Y K I I G R K N K \* V  
A F Y N E S Y G D + C Q E S F I L K K C F I S # F V E N I K E F A

misc_feature	27586	27591	gtattt, splice donor sequence
misc_feature	27862	27872	ttaacaatcg, splice branch and acceptor
misc_feature	27895	27900	gtaata, splice donor sequence
misc_feature	27972	27973	
misc_feature	28005	28006	
misc_feature	28404	28405	
misc_feature	28434	28435	gtaccc, splice donor sequence
misc_feature	28703	28715	ctgacaagtatag, splice branch and acceptor
misc_feature	28748	28753	gtaaat, splice donor sequence
misc_feature	28873	28885	ttaaccgataaaag, splice branch and acceptor
misc_feature	28938	28943	gtaagg, splice donor sequence

Artemis



Sanger

<http://www.sanger.ac.uk/science/tools/artemis>

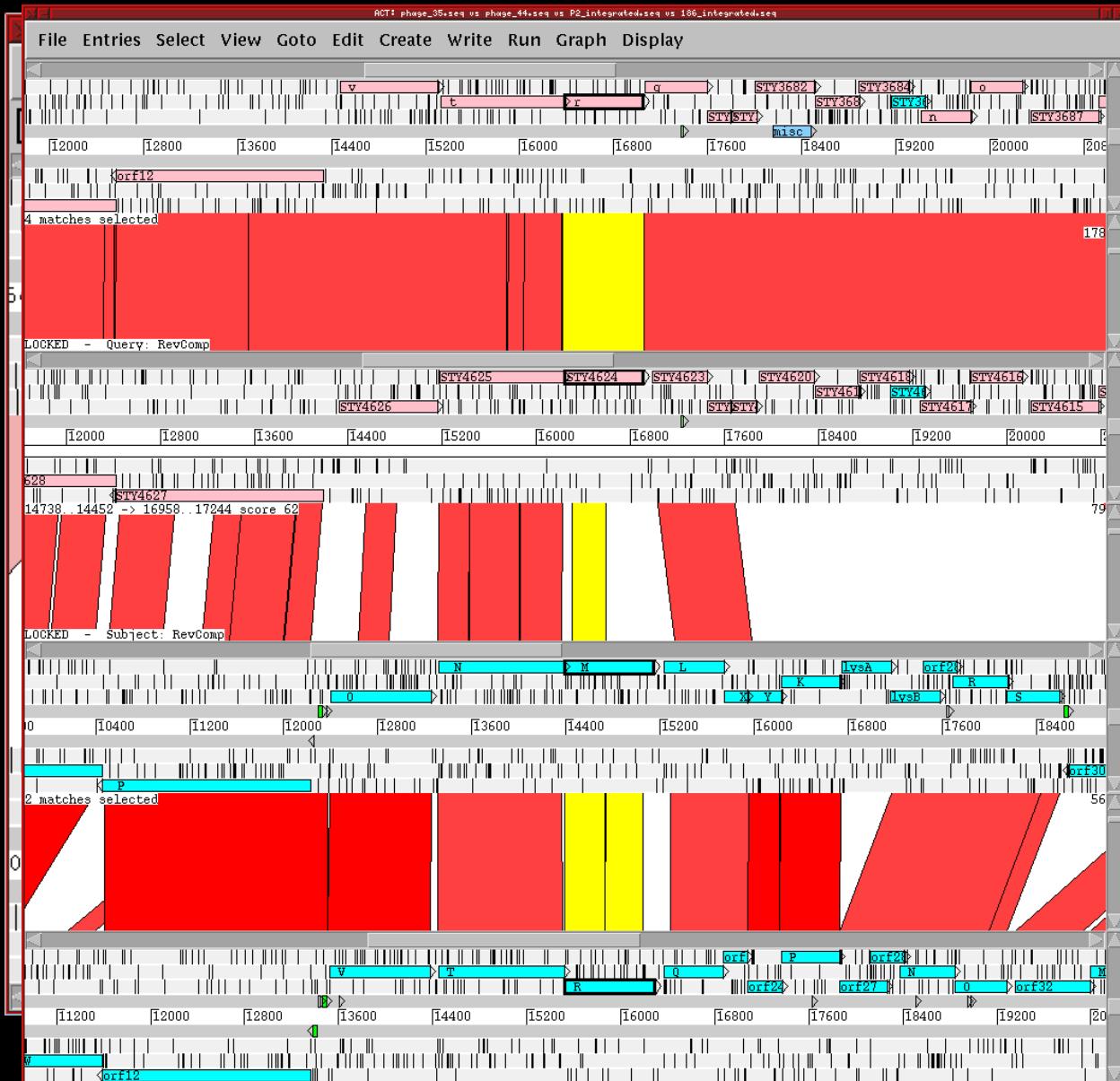
# Genome browser and annotation tool



Plug-in  
BamView

RNA-seq  
integration

# Comparative genomics for annotation



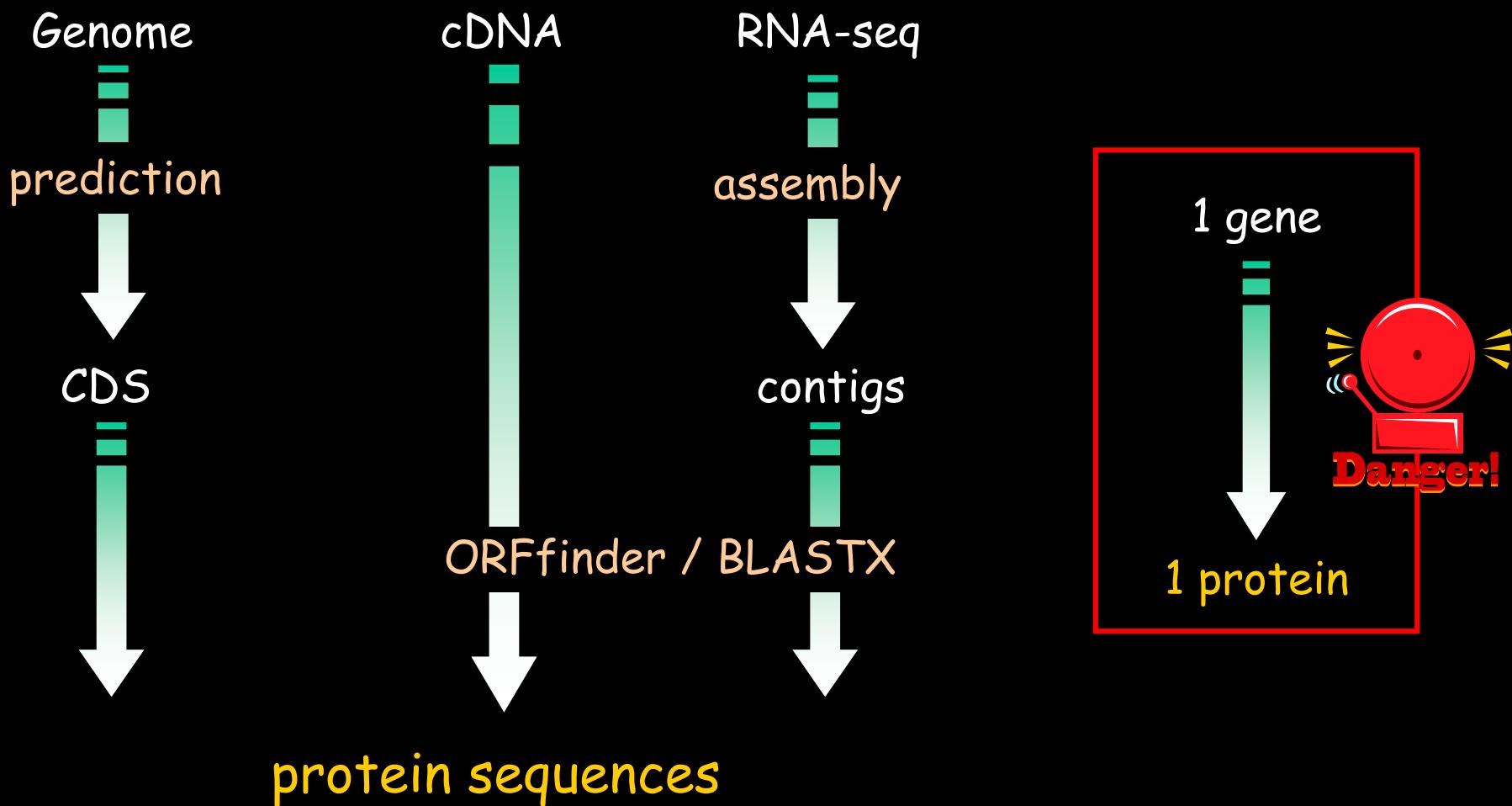
Plug-in  
ACT

---

# Genome Annotation

# Function

# Translation : from nucleotides to proteins



# Automatic functional annotation

## Prediction of protein function

- inference by similarity (BLAST)
- motif and domain detection

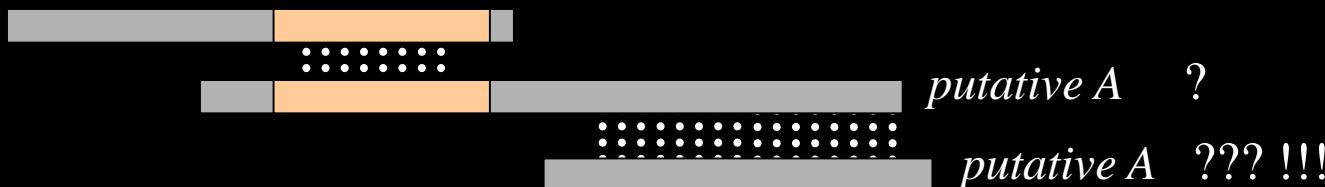
“ homologue to putative protein similar to hypothetical orphan - like ”

no hit	hypothetical	
hit RNAseq	unknown/expressed	◀ nomenclature
hit protein	putative/-like	



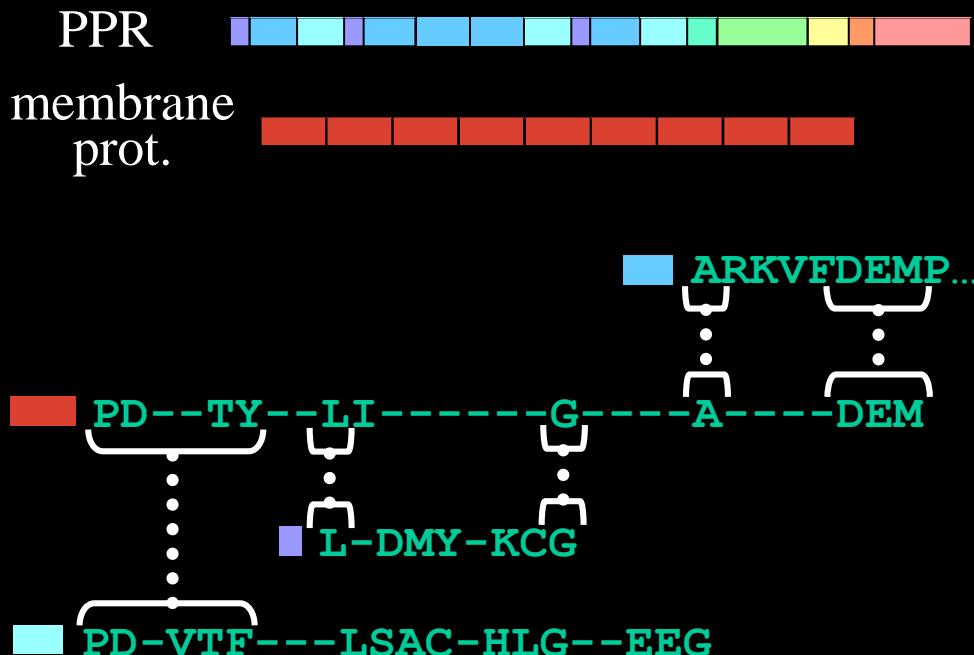
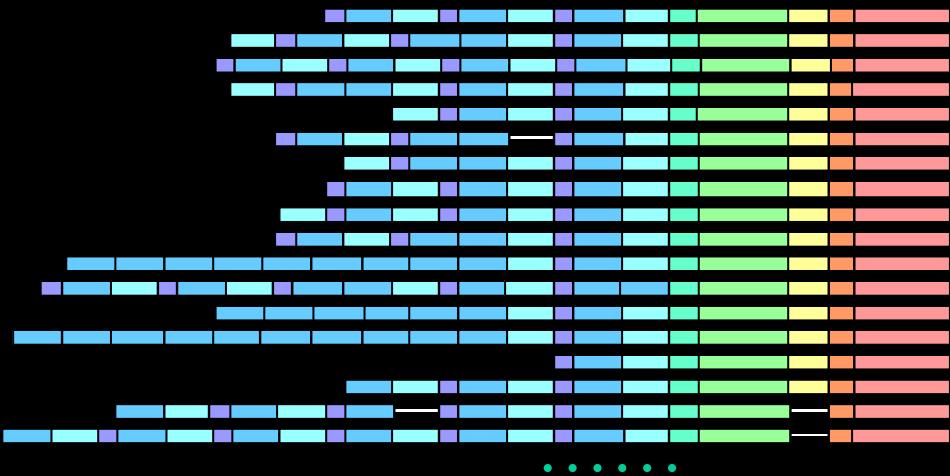
## Errors and databank contamination

*function A*



# PPR family

400 genes  
in Arabidopsis



## Annotations

- PPR protein
- selenium binding
- membrane associated



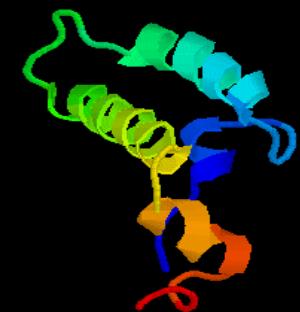
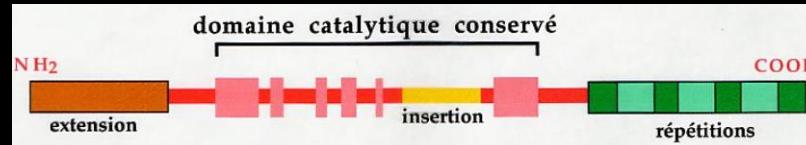
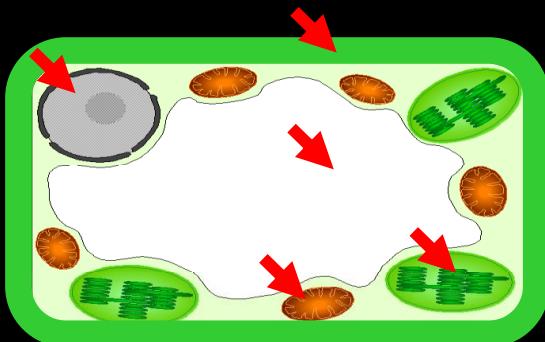
# Expert functional annotation

« Swiss-Prot method »

Bioinformatics predictions



function



- ✓ targeting and signal peptides for subcellular localization
- ✓ transmembrane segments
- ✓ motifs, domains and signatures
- ✓ post-translational maturation sites
- ✓ secondary structures and 3D models (PBIL)

# Motifs and domains

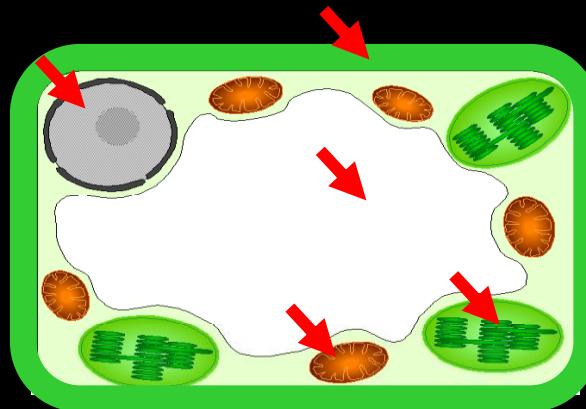
## Signal searching

TmPRED TmHMM

PSORT

PREDOTAR

TargetP MitoP ChloroP...



## Known motif searching

Prosite (SIB-expertised)

1 500 entries

Pfam-AB (profiles hmm-GeneWise)

12 000 entries

Prints (footprints - signatures)

2 000 entries

Prodom (automatic clustering)

574 000 entries

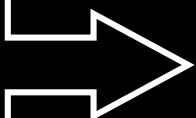
Panther (evolution-expertised)

6 000 entries

# Motifs and domains

Looking for known motifs

UNIPROT  
PFAM  
PRINTS  
PROSITE  
PRODOM  
SMART  
SUPERFAMILY  
TIGRFAM  
PIRSF  
GENE3D  
PANTHER

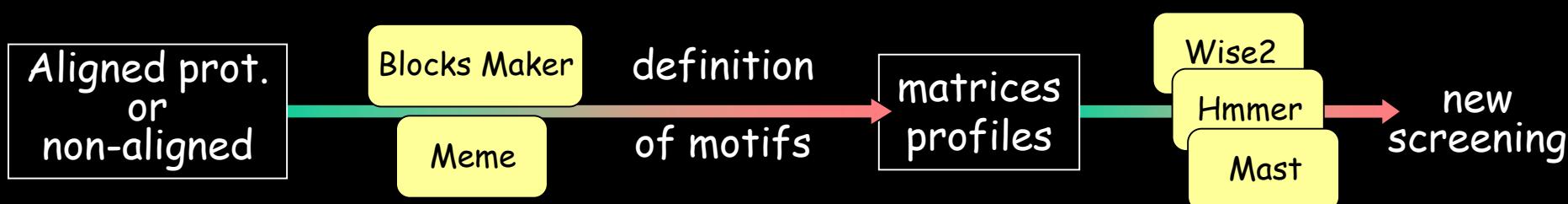


## InterPro

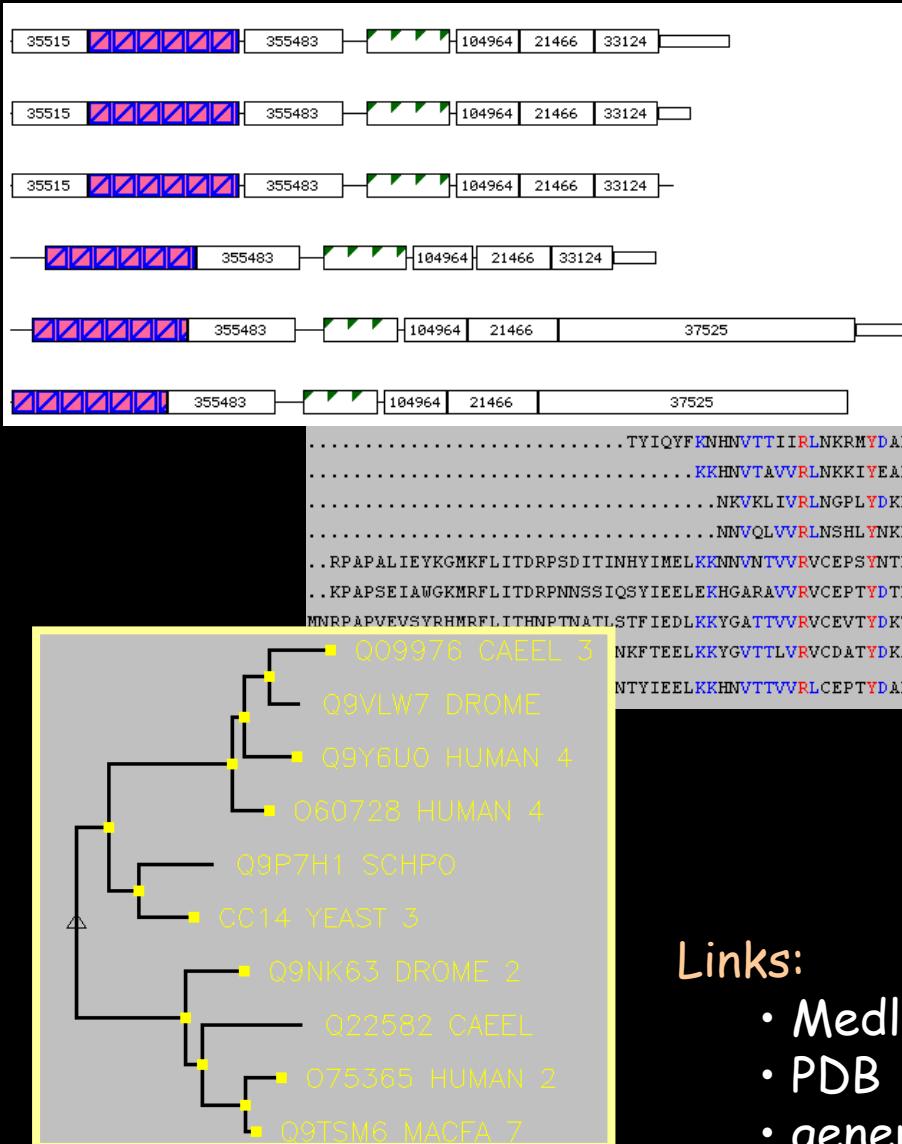
<http://www.ebi.ac.uk/interpro/index.html>  
Nucleic Acids Research vol 29(1):37-40

Release 39  
23 800 entries  
24 10<sup>6</sup> proteins  
(80% trEMBL)

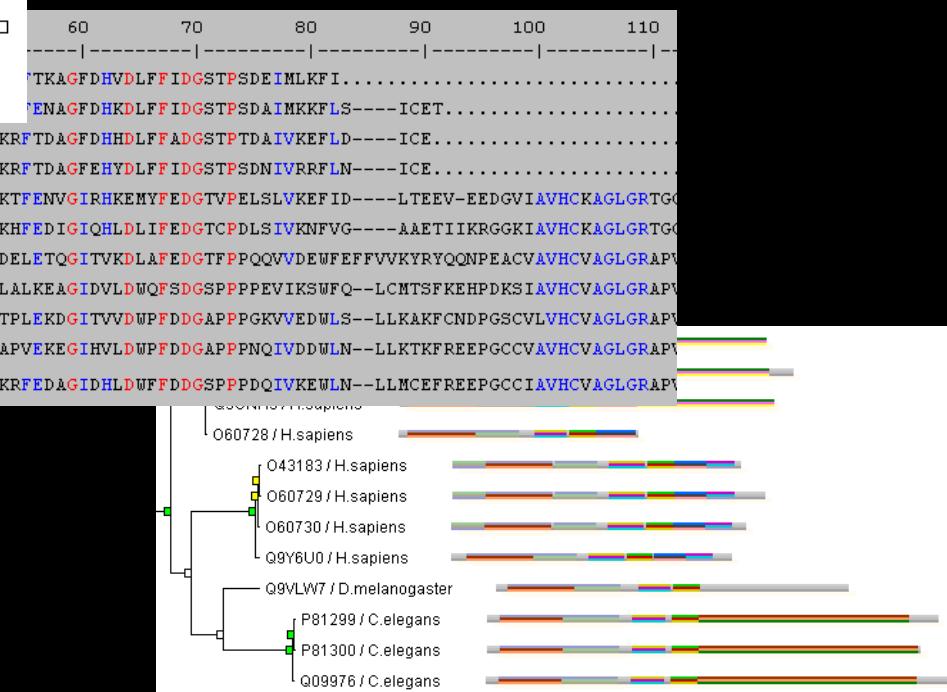
Definition of new motifs



# INTERPRO / INTERPROSCAN : the reference



- phylogeny
- cross-references
- biochemical function
- examples



## Links:

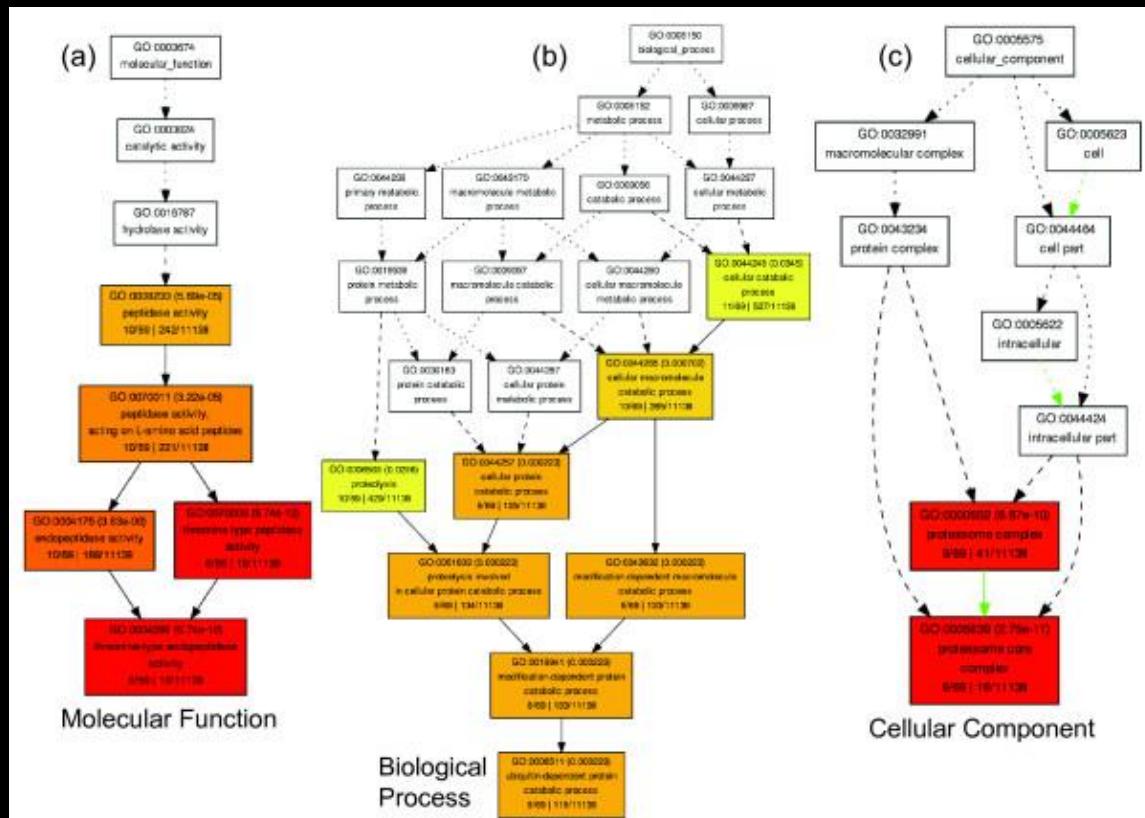
- Medline
- PDB
- generic databanks
- ...

# Functional classification of genes



## Gene Ontology Consortium

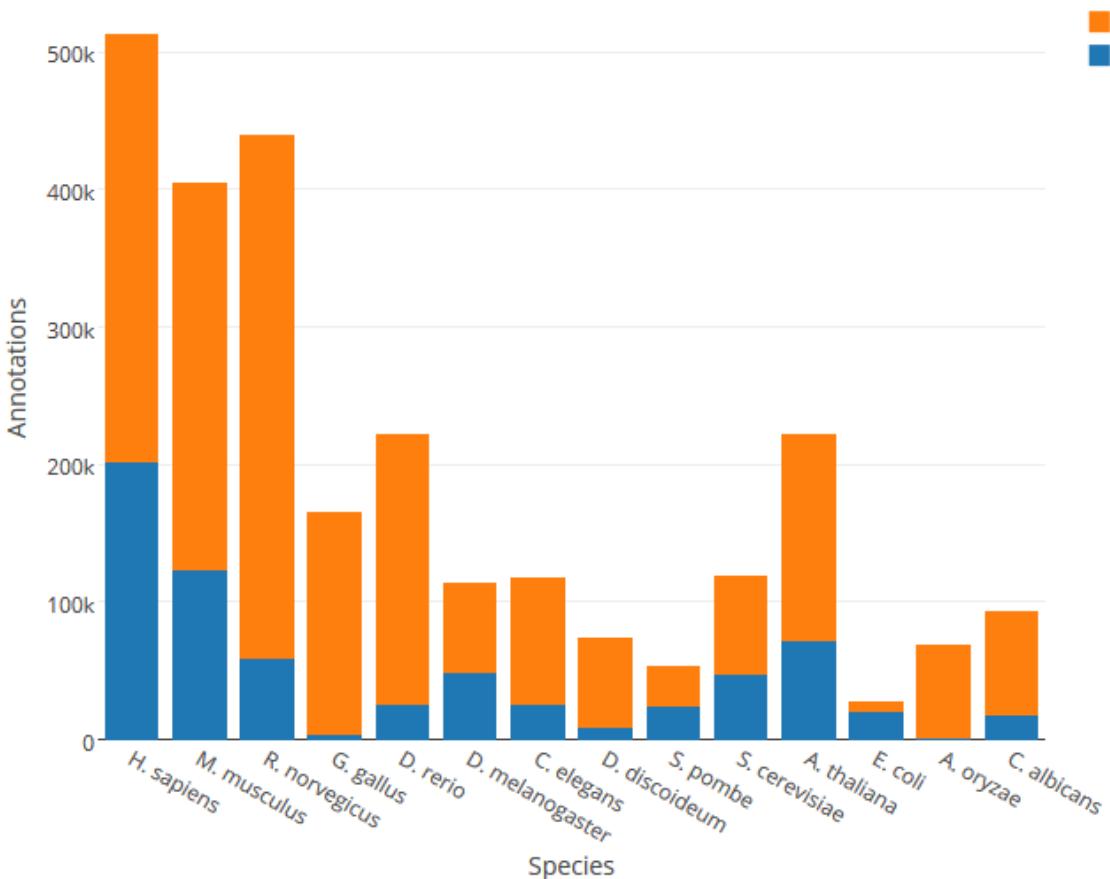
- molecular function (transcription factor, protease...)
- biological process (mitosis, sugar metabolism...)
- cellular component (nucleus, ribosome...)



BIOLOGICAL PROCESS
Cellular component organization and biogenesis
Nucleobase-containing compound metabolic process
Multicellular organismal development
Response to stress
Transport
Cell communication/Cell-cell signaling
Protein metabolic process
Reproduction
Anatomical structure morphogenesis
Signal transduction
Cell differentiation
Lipid metabolic process
Cell cycle
Cellular protein modification process
Carbohydrate metabolic process
DNA metabolic process
Response to external stimulus
Post-embryonic development
Death
Growth/Cell growth
Response to endogenous stimulus
Embryo development
Nucleic acid binding
Cellular homeostasis
Response to biotic stimulus
Response to extracellular stimulus
Translation
Regulation of gene expression, epigenetic
Secondary metabolic process
Generation of precursor metabolites and energy
Flower development
Photosynthesis
Pollination/Pollen-pistil interaction

# Functional classification of genes

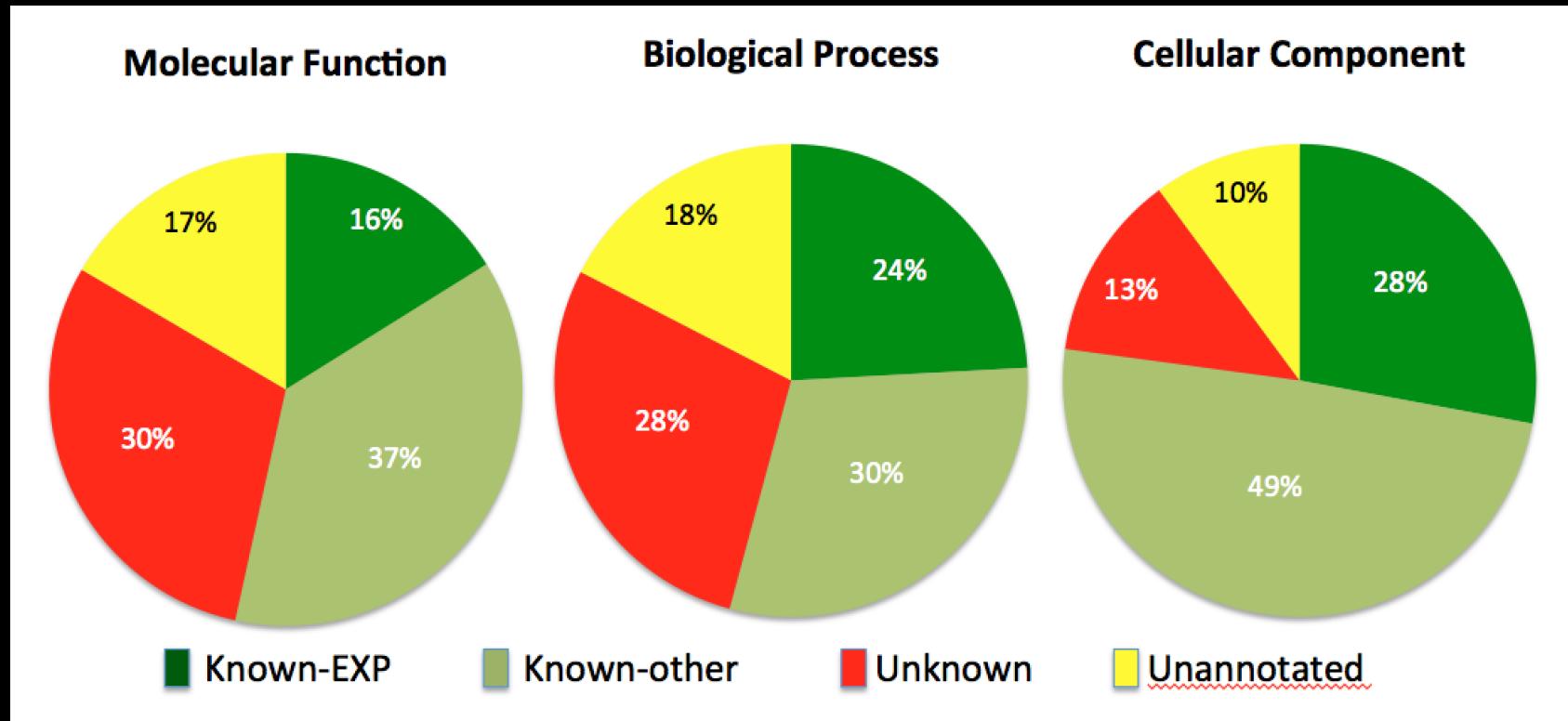
Experimental annotations by species



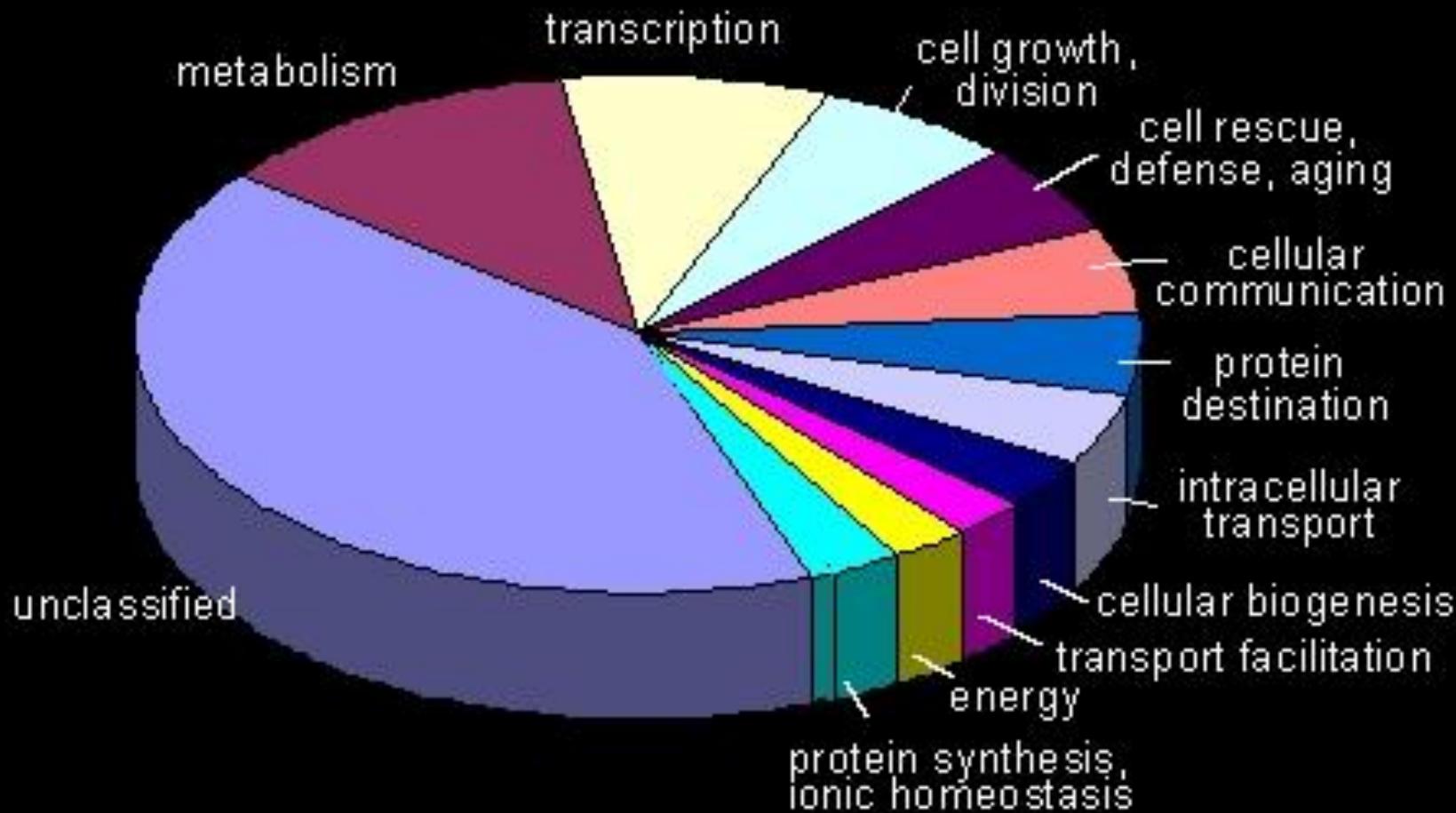
Inferred from:

- Experiment (EXP)
- Direct Assay (IDA)
- Physical Interaction (IPI)
- Mutant Phenotype (IMP)
- Genetic Interaction (IGI)
- Expression Pattern (IEP)
- Sequence or structural Similarity (ISS)
- Sequence Orthology (ISO)
- Sequence Alignment (ISA)
- Sequence Model (ISM)
- Genomic Context (IGC)
- Biological aspect of Ancestor (IBA)
- Biological aspect of Descendant (IBD)
- Key Residues (IKR)
- Rapid Divergence (IRD)
- Reviewed Computational Analysis (RCA)
- Curator (IC)
- Electronic Annotation (IEA)

# Functional annotation of the *Arabidopsis* genome



# Functional annotation of a whole genome



# UniProt & Swiss-Prot

UniProt KnowledgeBase



TrEMBL

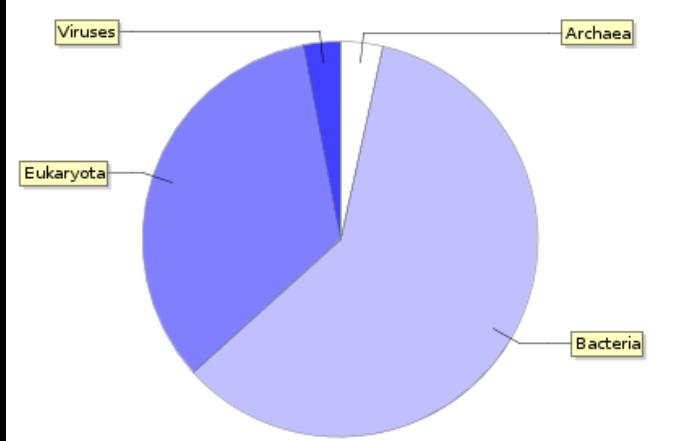
Release 2018\_07  
120 10<sup>6</sup> entries

Swiss-Prot

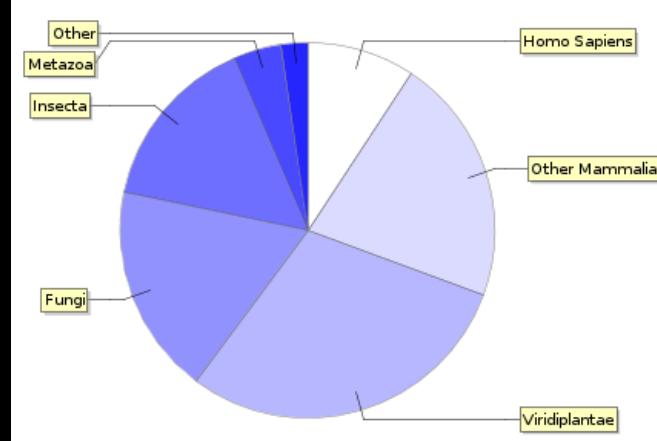
Release 2018\_07  
558 000 entries (2%)

- 50 annotators (SIB) + experts
- Monthly updated
- Curated predictions /validation
- Data mining

Swiss-Prot entries per taxonomic group



Swiss-Prot entries in Eukaryota

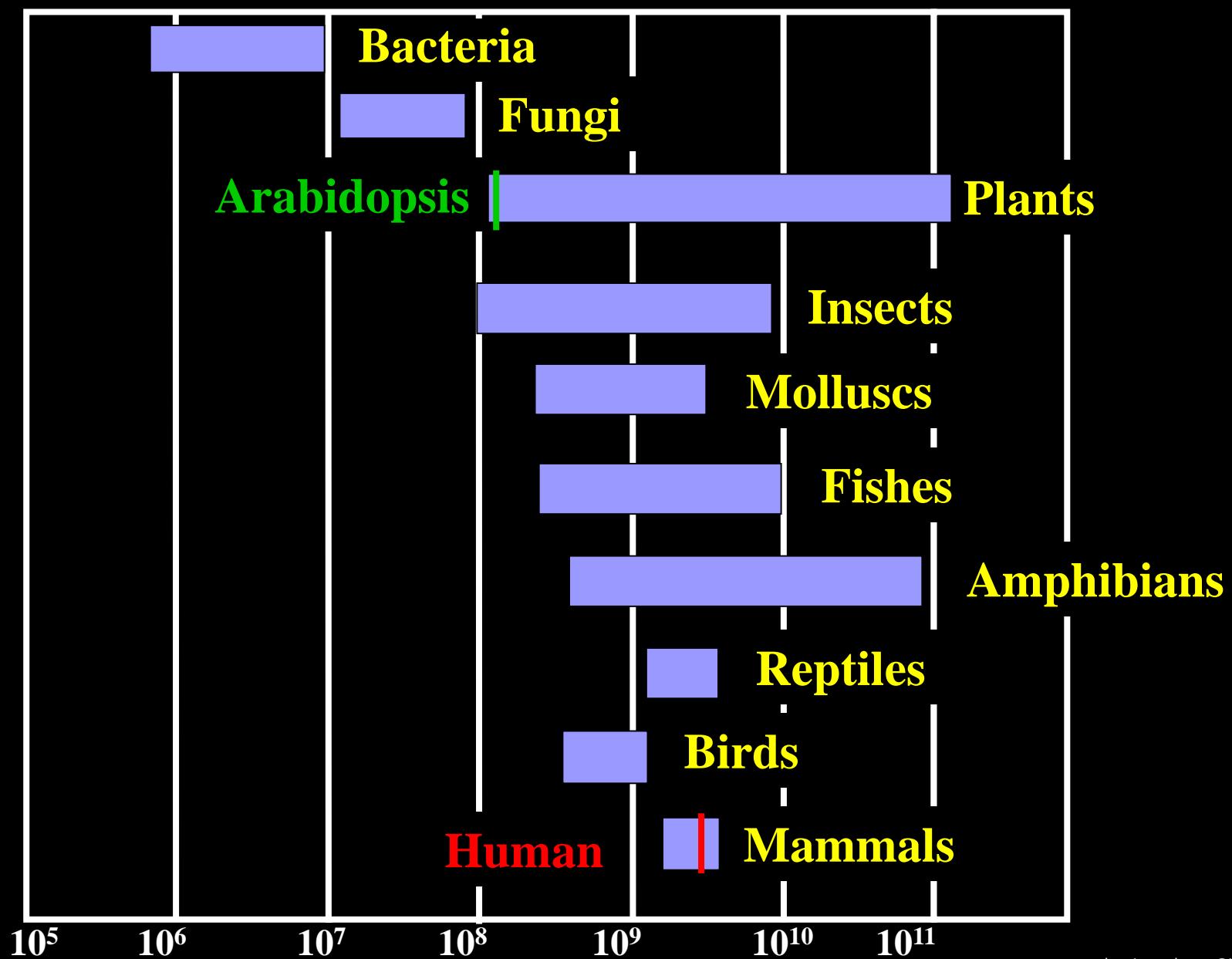


---

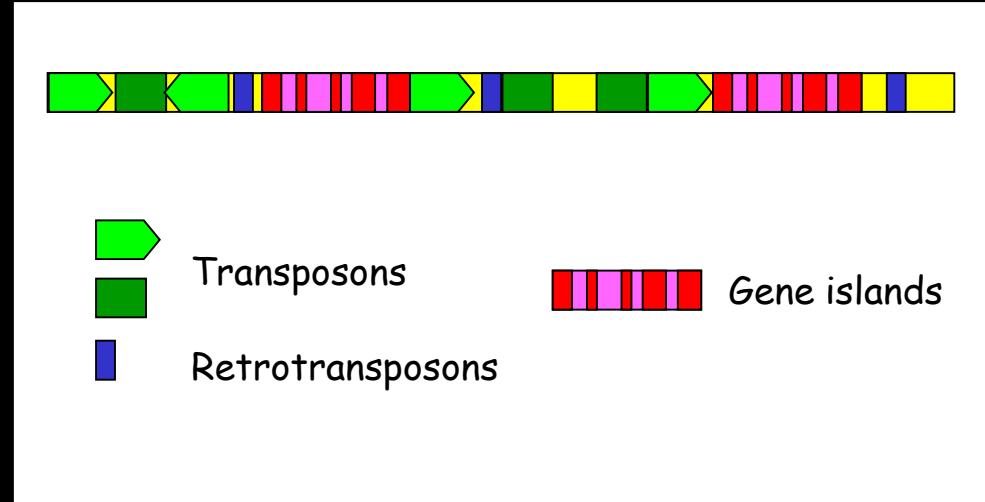
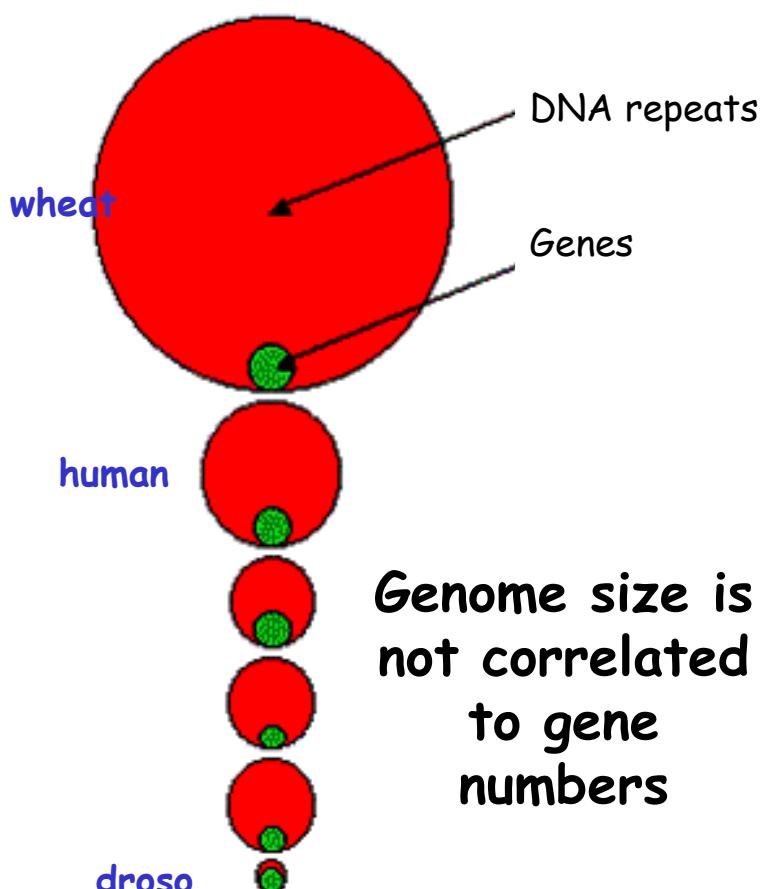
Genome Annotation

Evolution

# Genome size



# Gene Space

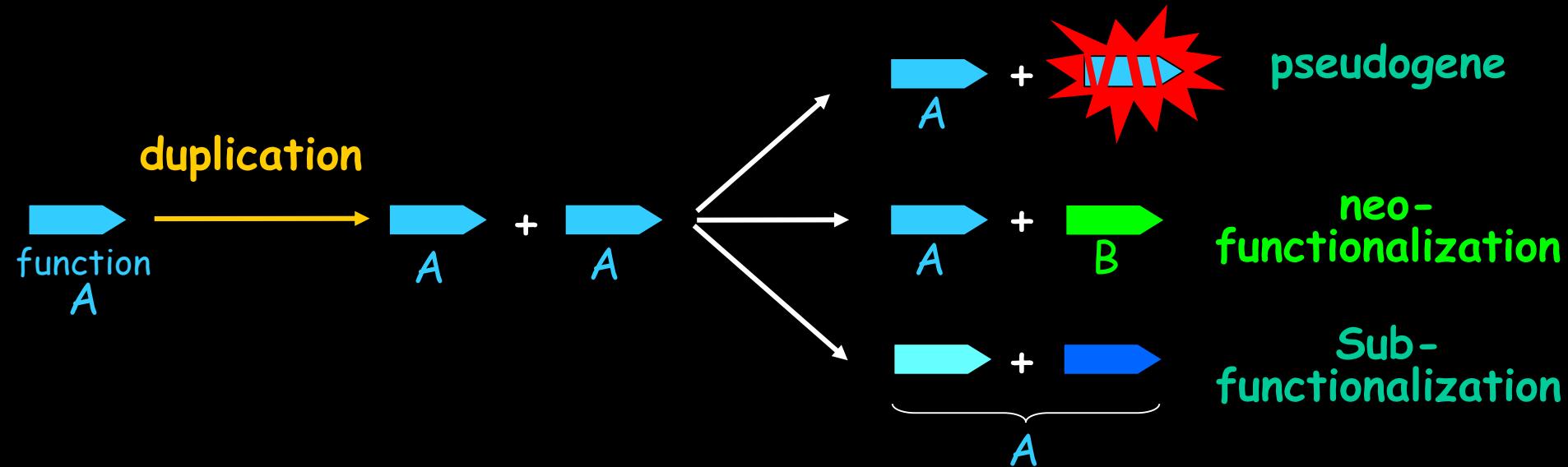


Genome size is the results of opposite processes :

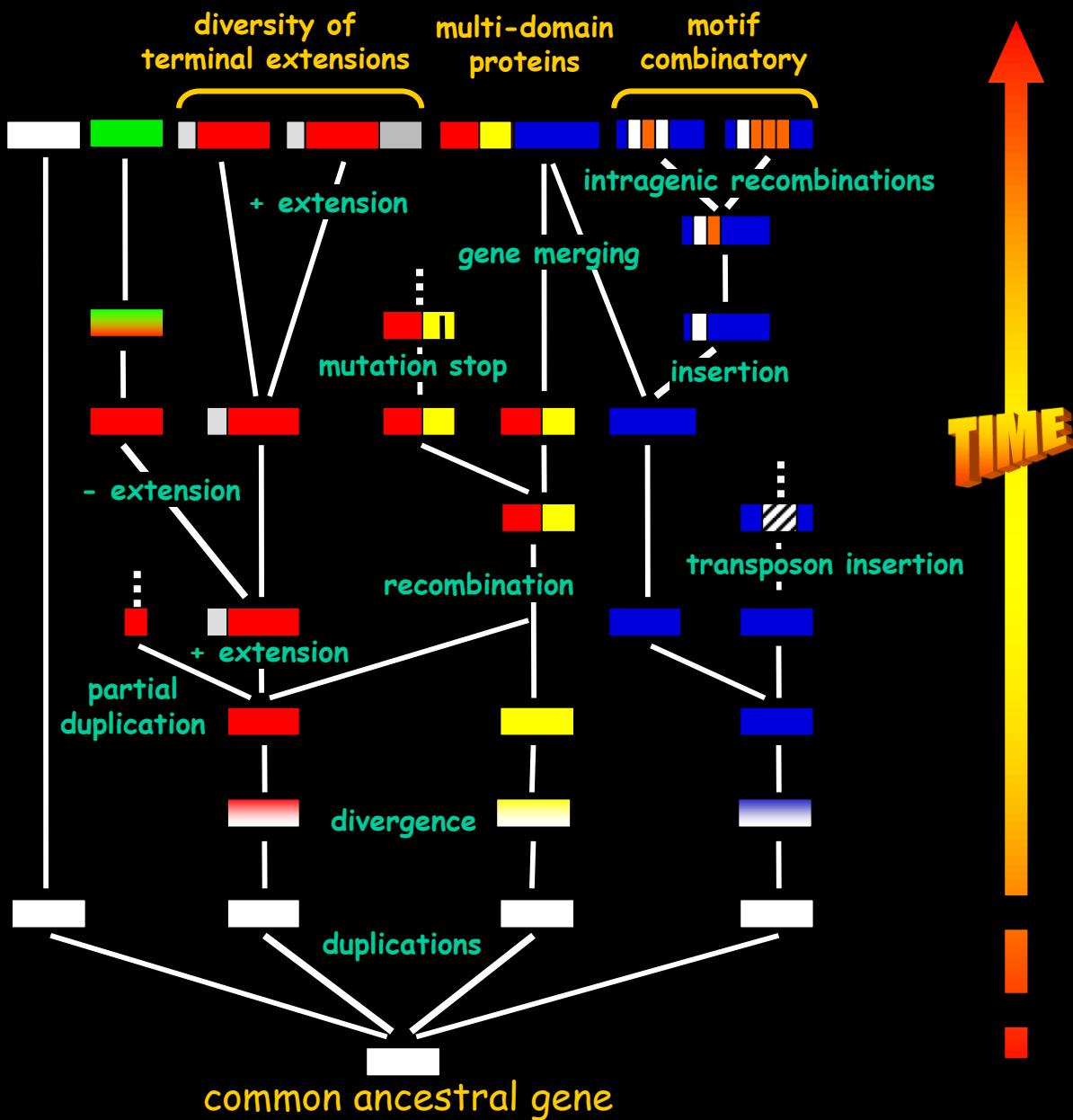
- Retrotransposon invasion by multiple insertions
- Duplications (at different scales)
- Deletions by recombination

# Duplication and divergence

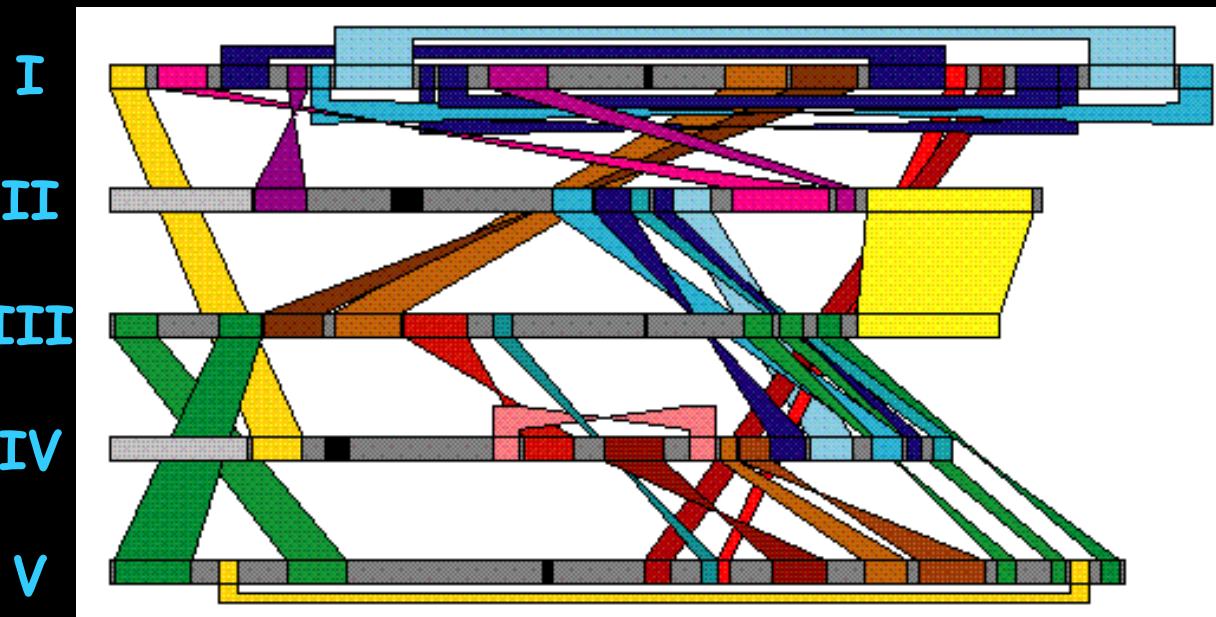
'Evolution by gene duplication'  
Ohno, 1970



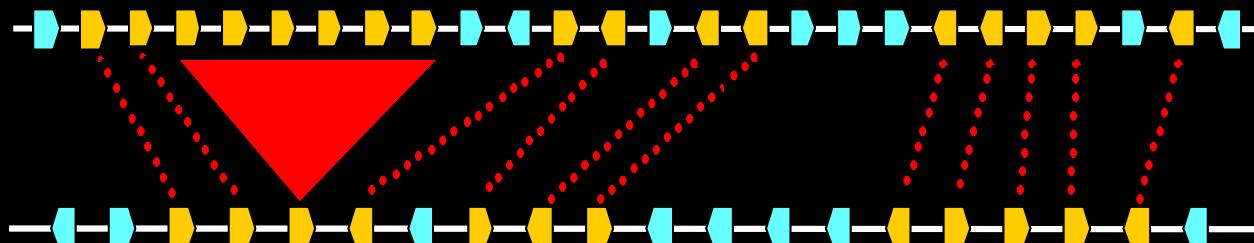
# Functionalization



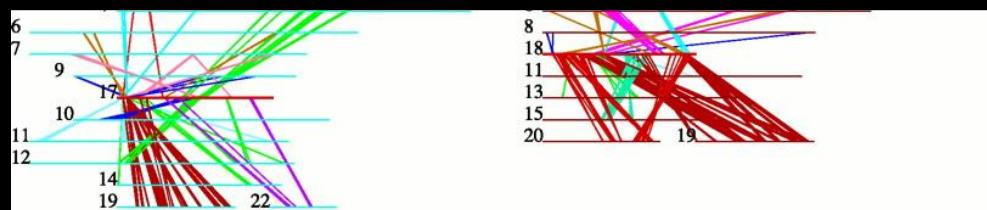
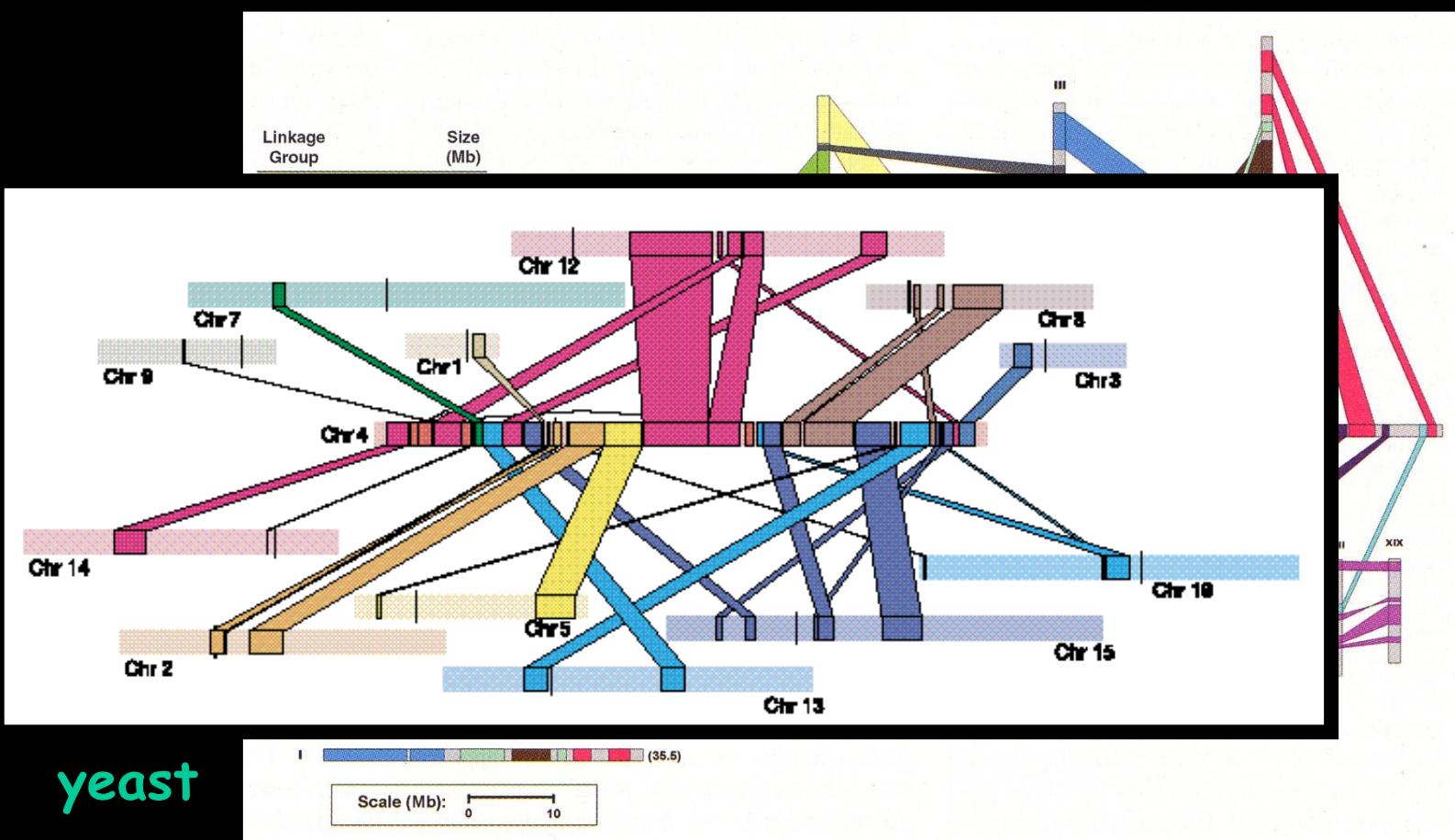
# Dynamic of genomes



80% of the  
Arabidopsis  
genome is  
duplicated



# Whole Genome Duplication



# Homology

---

Etymology : A homologous trait is any characteristic of organisms that is derived from a common ancestor, a common origin. The homology is a similarity acquired, inherited from a common ancestor.

The significant similarities between two sequences are considered as representative of their homology: They have a common origin, they derived from the same ancestral sequence.

Similarity (quantitative)  Homology (qualitative)

# Definitions

---

## HOMOLOG

### PARALOG

Homologous genes are paralogs if they were separated by a gene duplication event.

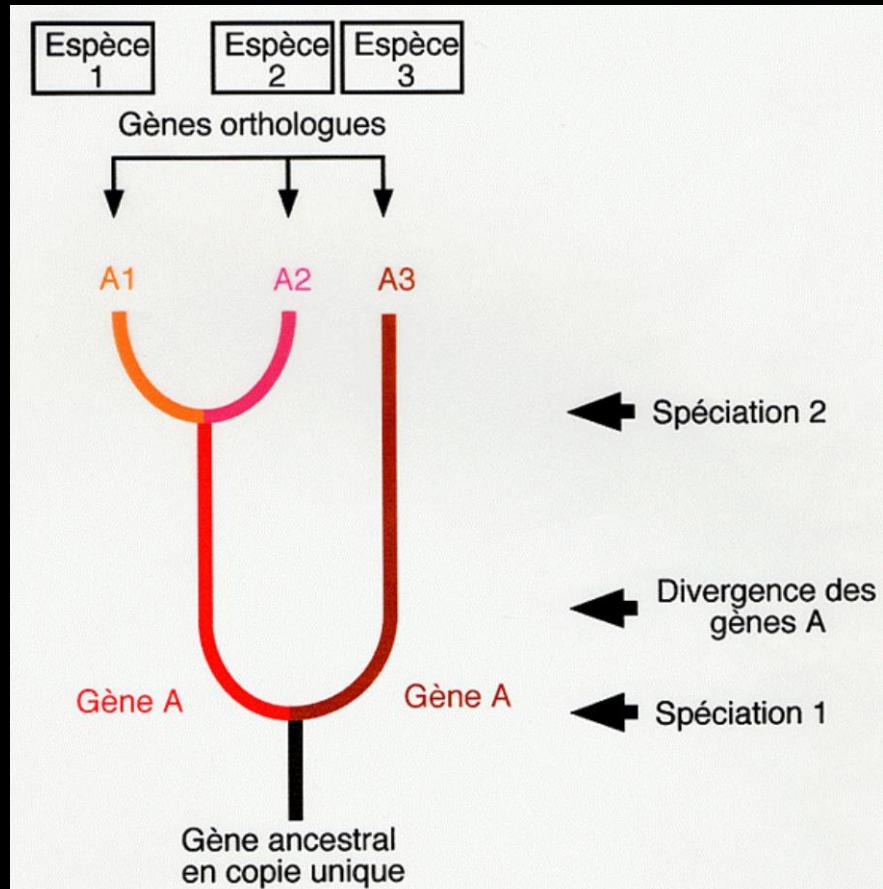
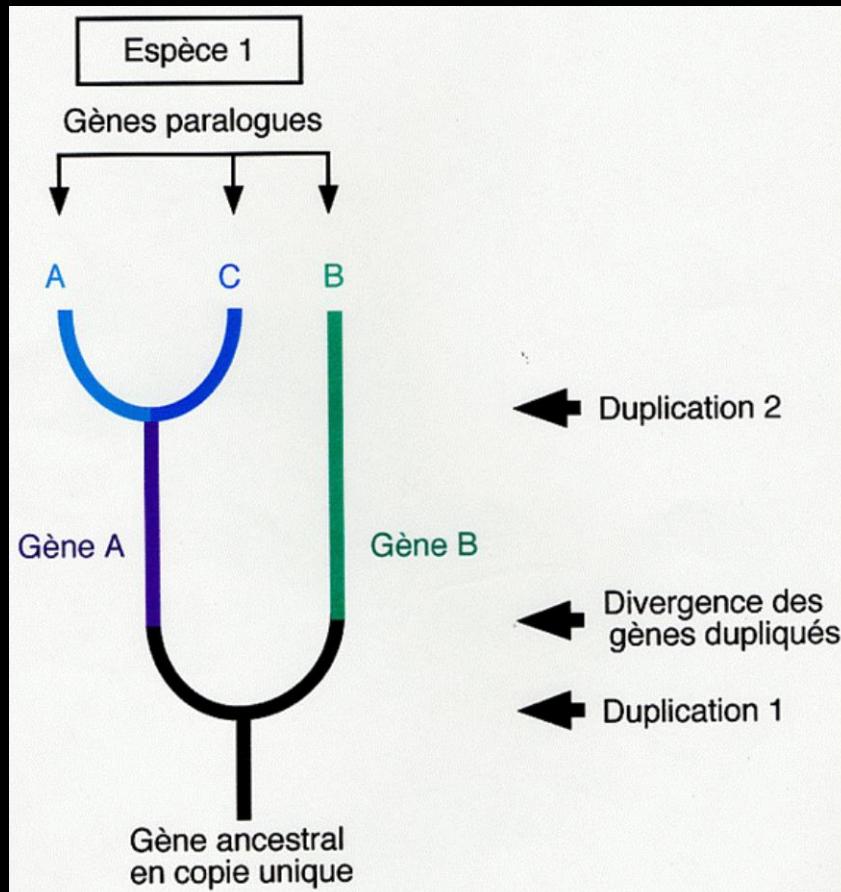
### ORTHOLOG

Homologous genes are orthologs if they were separated by a speciation event: when a species diverges into two separate species, the divergent copies in the resulting species are orthologs.

## GENE FAMILY

Group of homologous genes. The genes belonging to the same family have therefore a common ancestor gene and share significant similarities. They encode proteins having a common biochemical function but can perform different biological functions.

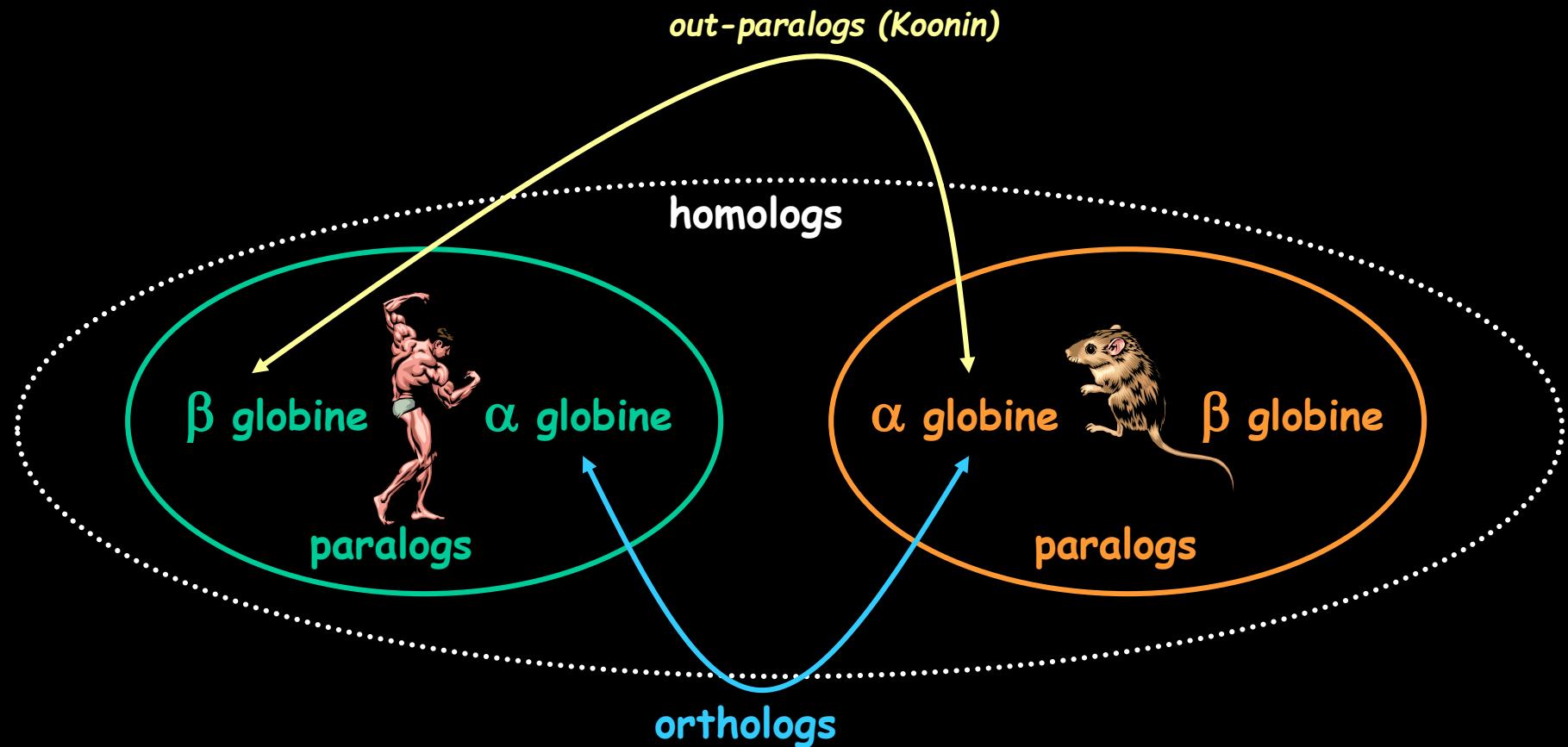
# Paralogs and orthologs



# Concretely...

The genes AtMYB29 and OsMYB43 have 53% of similarity in their coding region.

- similarity
- identity
- homology

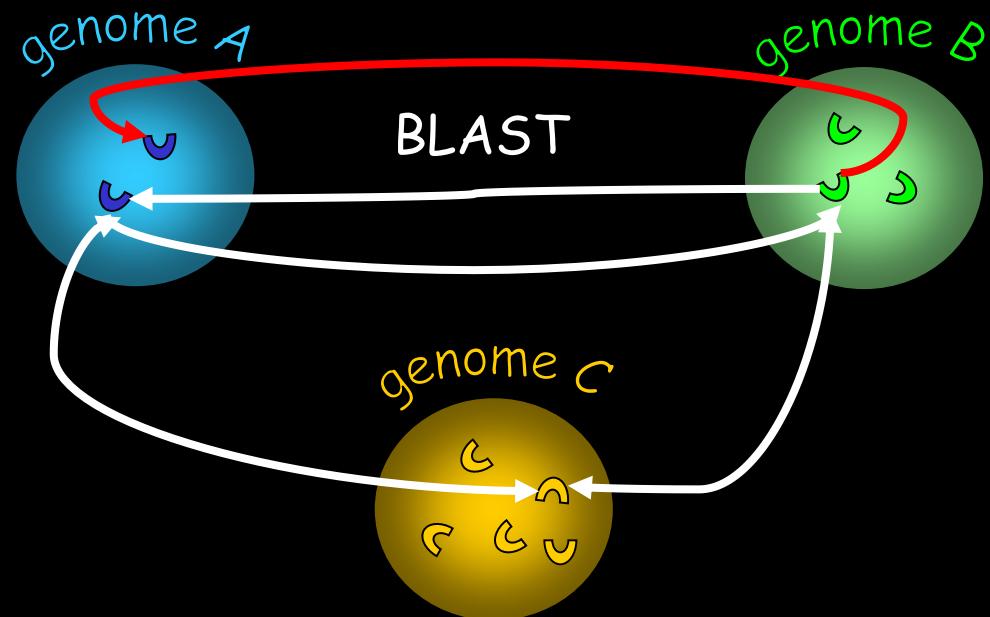


# The ortholog quest ('functional orthologs')

Diversify the approaches... collect the clues

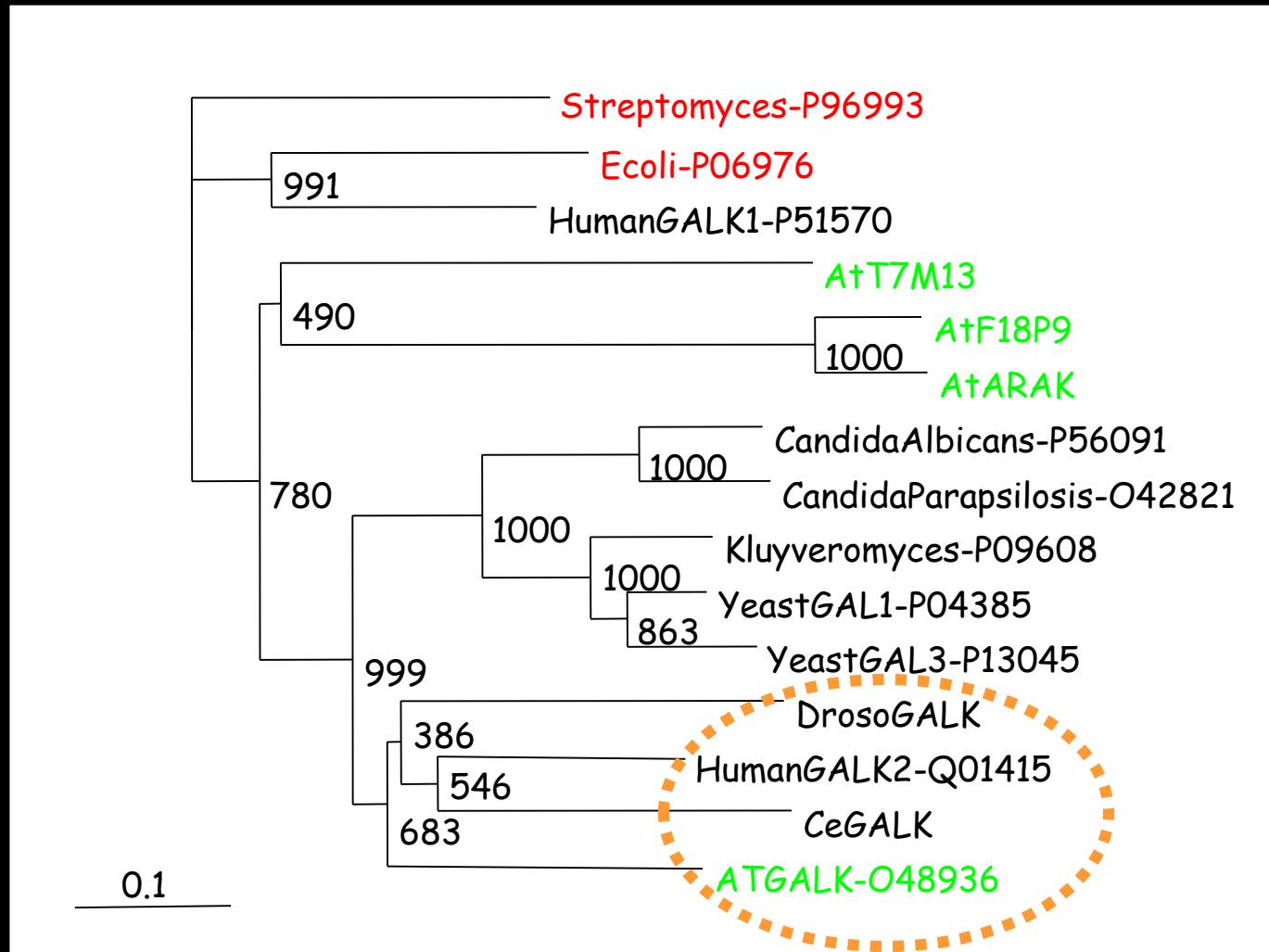
- Cross the genomes:  
'BDBH'

bidirectional best hit  
multidirectional reciprocal best blast hit

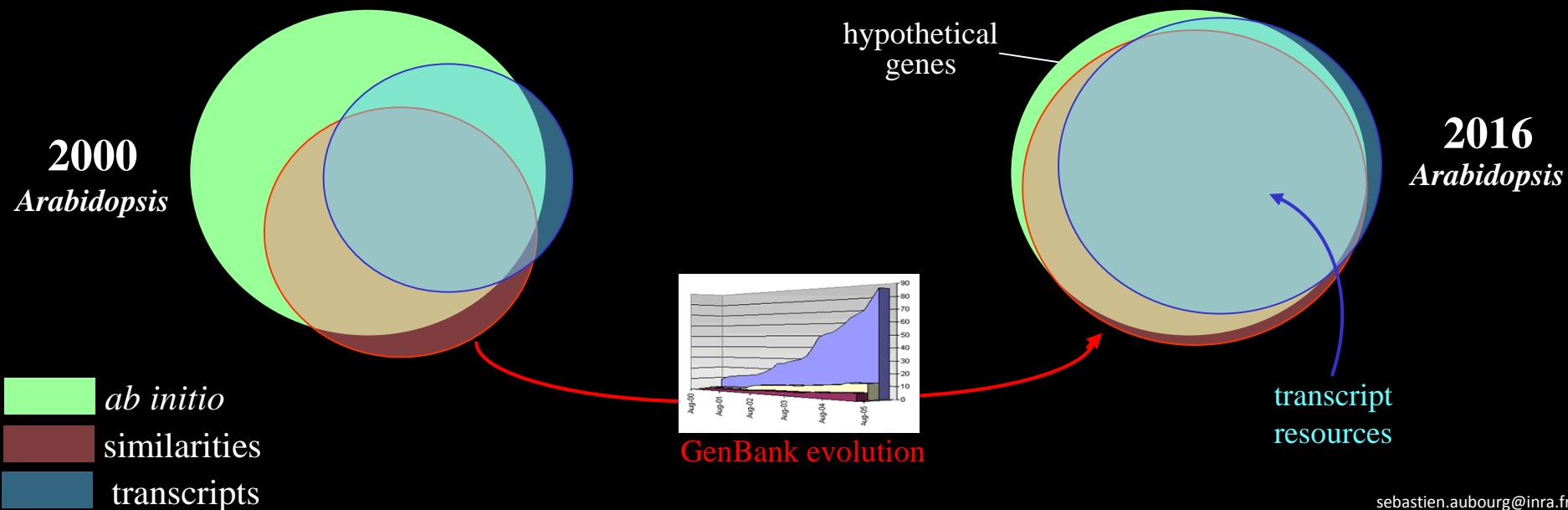
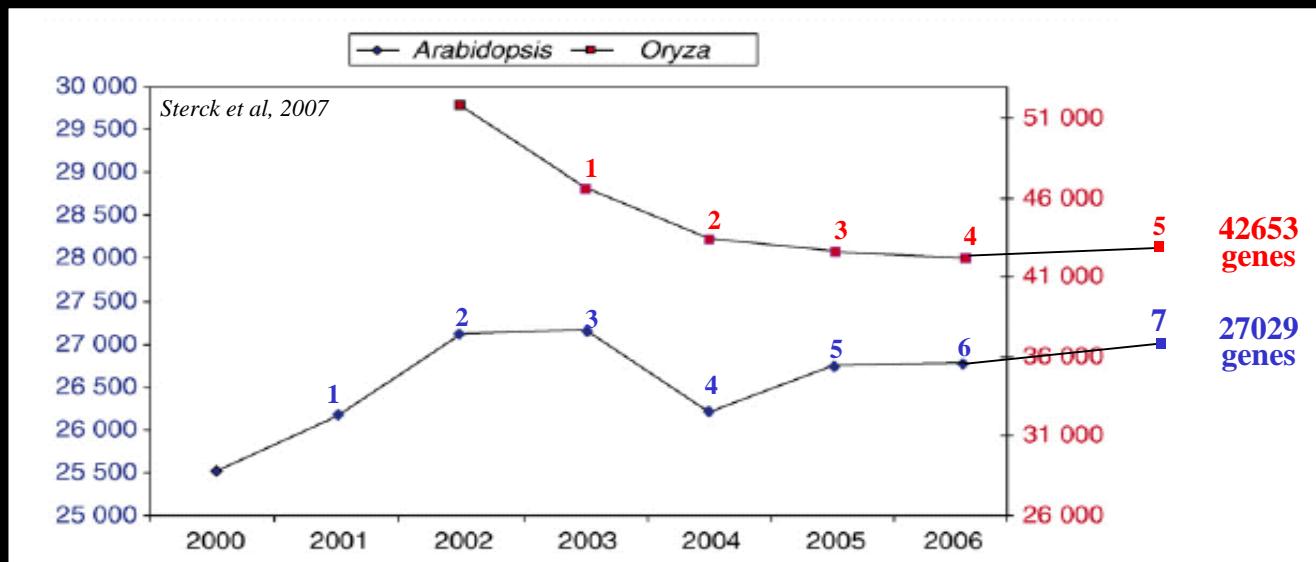


- Look for the presence of common insertion or extension
- Compare intron-exon structure of genes
- Find protein partners (Rosetta Stone)
- Exploit the genomic environment (synteny)
- Explore promoter regions (promoter shadowing)
- Exploit transcriptomic or proteomic data (expression profiles)

# Phylogeny and orthology



# Annotation is a dynamical process



# Toward a relational annotation

---

- Comparative genomics (gene context)  
Phylogenetic profiles ► functional collaboration
- miRNA prediction and their target(s)  
Relationships between TFBS/TF (ChIP-) ► regulation
- Transcriptome  
Proteome ► co-expression
- Interactome (systematic Y2H)  
Definition of interologs ► direct interaction
- Text-mining (co-citation) ► functional links

► gene/protein networks  
Integrative biology

# From sequences to functions

