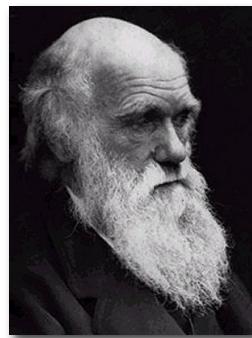
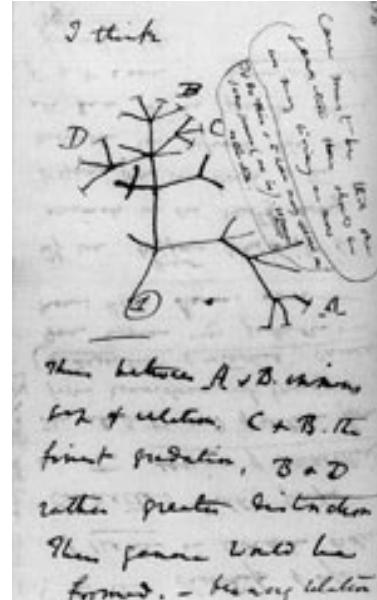


Evolution des génomes eucaryotes:

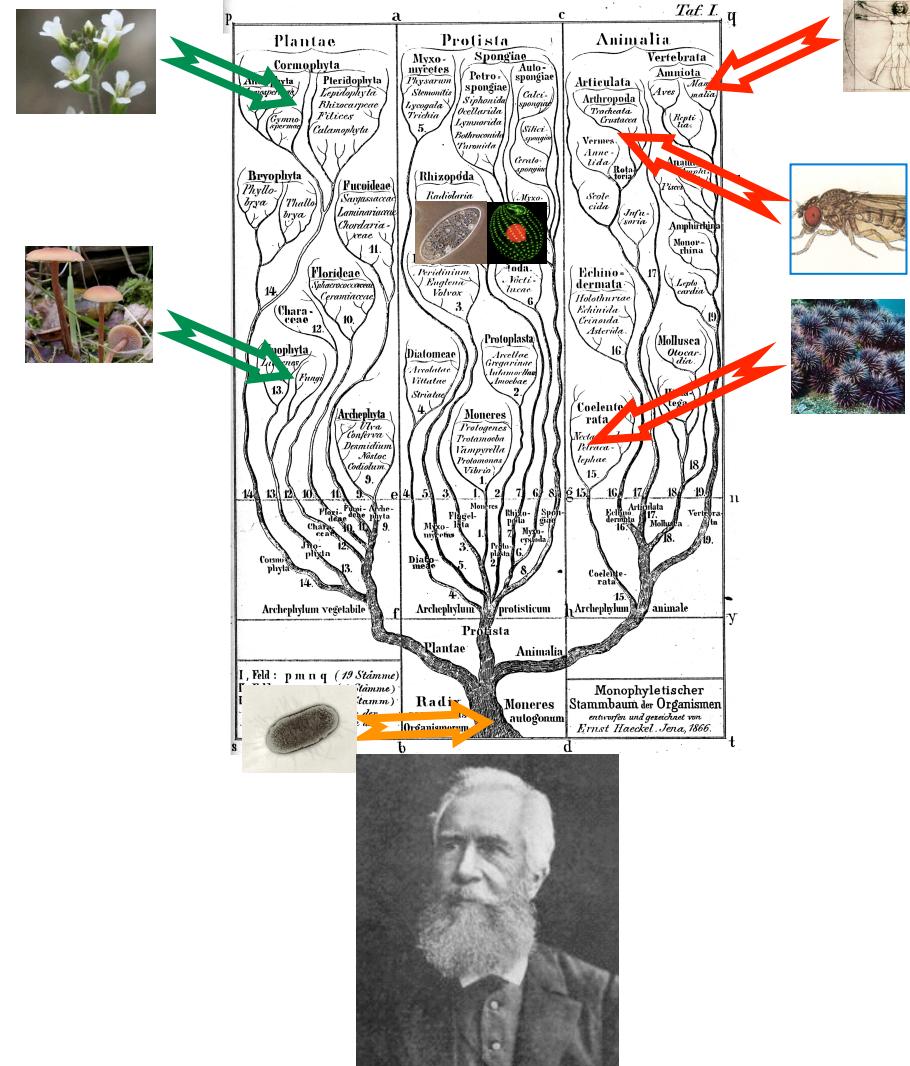
exemple des levures et autres unicellulaires

Cours BIM-PHYG
4 décembre 2014
Bernard Dujon



Charles Darwin
Reproduction avec
changement graduel

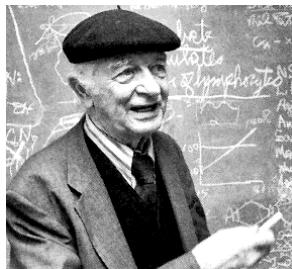
- 1- Les organismes vivants n'évoluent pas, ils se reproduisent. Ce sont les populations qui évoluent.
- 2- La génomique moderne est en train de bouleverser nos conception sur l'évolution.



Ernst Haeckel, 1866
Phylogénie / ontogénie

1- Les séquences des protéines homologues divergent

Q96J94/555-847	IVVCLLSS.NRKDK.YDAIK KYLCTDCPTPSOCV VARTLGKQ.....QT.....VMAIAT KIALQMNC KMGGE...LWRVDI
Q23415/40-350	KFAF VITD .DSITH.LHKKY KALEQKSMMVIQDMKISKAN SV.....VK.....DGKR LTLENVINKTNM KLGGLN..YTVSDI
Q74957/500-799	YLFF FILDK .NSPEP.YGSIK RVCNTMLGVP SQCAIS KHILQSKPQYC ANLGMKINVKV GGINC.SLIPK
Q62275/594-924	TFVFI ITD .DSITT.LHQRY KMIEKDT KMIVQDMKLSKALSV.....INAGKR LTLENVINKTNV KLGGSN..YVFVDI
Q58717/426-699	CFALI IGKEKYKDNDYYE ILKKQLFDLKI I SONILWENWRKD.....DK.....GYMTNNLLIQIM GKLGK IK...YFILD
Q9VKM1/538-829	LILCL VPN .DNAER.YSSIKR GYVDR AVPTQVVTLKTTKNR.....SL.....MSIAT KIAIQLNCKL GYT...PWMIE
Q21691/673-1001	TIV FGIIA .EKRPD.MHDIL KYFEEKLGQQTIQ ISSETADKF.....MR.....DHGGK Q TDN VIRKLNPKCGGTNF LIDVPE
Q16386/548-847	QLIM FITK ..SMNN.YHTEIK CLEQEFDLLTODIFETAVKLAQ.....QONTR KNI YKTNM KLGGLN ..YELRS
Q9XGW1/625-946	LLLAI LPD .NNGSL.YGDLKRICETELGLISQC CLTKHVFKISKOYLANV SLKINVKM GGRNT.VLVDA
Q17567/397-708	MLVV MLAD .DNKTR.YDSLKKYLC VECP IPNQC VNLRTL AGK.....SKDGG E NKNLGSIV LKV LQ OMIC KTGGA...LWKVN
Q19645/674-996	PFVL FISD ..DVPN.IHECL KFEERM SDIPTQ HV LLKNV KKM RDNIEKK SGGRR RAYDLTLDNIV VMK ANI KCGGLN ..YT.ADI
Q09249/650-977	DILVG IAR .EKKPD.VHDIL KYFEESIGL QTQ QLCQQT VDKM.....MG....GQGG Q RTD NV MRKF FN LKG GGT NFFV EIPN
Q21495/52-336	GIVLPTPRIFFRD GQETSLNNQSFRNPT.....DFAQTG FFV DAK QQLG GLN..YVVNS
Q28951/110-406	GIML VLP E.YNTPL.YYKL KSYLINS ..IPSQ FMRYD ILSNR.....NL.....TFYVD NLL VO FV SKLGGK..PWILN
Q16720/566-867	LIVV VLPG ..KTPI.YAEV KRVGDT VLGI ATOCVQAKNAIRTTPQ TLSNLC IKMN V KLGGVNS.ILLPN
Q76922/555-852	IVMV VVMRS .PNEEK.YSCI KRTCVD RPVPSQVVT TLKVI APR.....QQ....KPTGLMSIA T VVI QMN AKLMGA...PWQVW
Q9ZVD5/577-885	FILC VLPDKKN SDL.YGPWKK KLTEFGIV TOCMAP TRQP ND.....QYLTNLL KINA KLGG LN MSLVEI
Q48771/542-860	FILC I LPERK TS DI.YGPWKK ICLTEEGI HTOCICPIKI.....SDQYLTNV L LIKINS KLGG INS.LLGIE
Q67434/419-694	LVIV FLEEYPKVDP .YKSFL LYDFVKRELLKKM IPSQVILNR.....TL...KNE NLFV LLNV AEQVLA KTGNIP..YKLKE
P34681/660-966	CII VVLQS .KNSDI.YMTV KEQSDIVHGIM SQCVLM MKN NSRP.....TPATCANIV KL N MKG GIN..SRIVAI
Q02095/574-878	QLLFFVV K ..SRYN.YH QQI KALEQ KYDVLT QEIRAETA EKVFR.....QPQ TRLN I I NKTNM KLGG LN..YAIGS

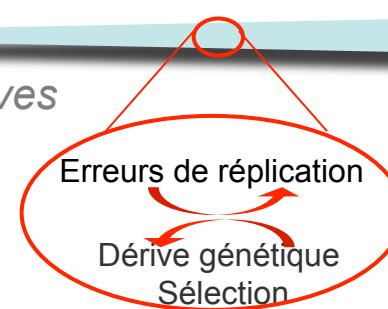


Linus Pauling

Horloges moléculaires

Phylogénies moléculaires

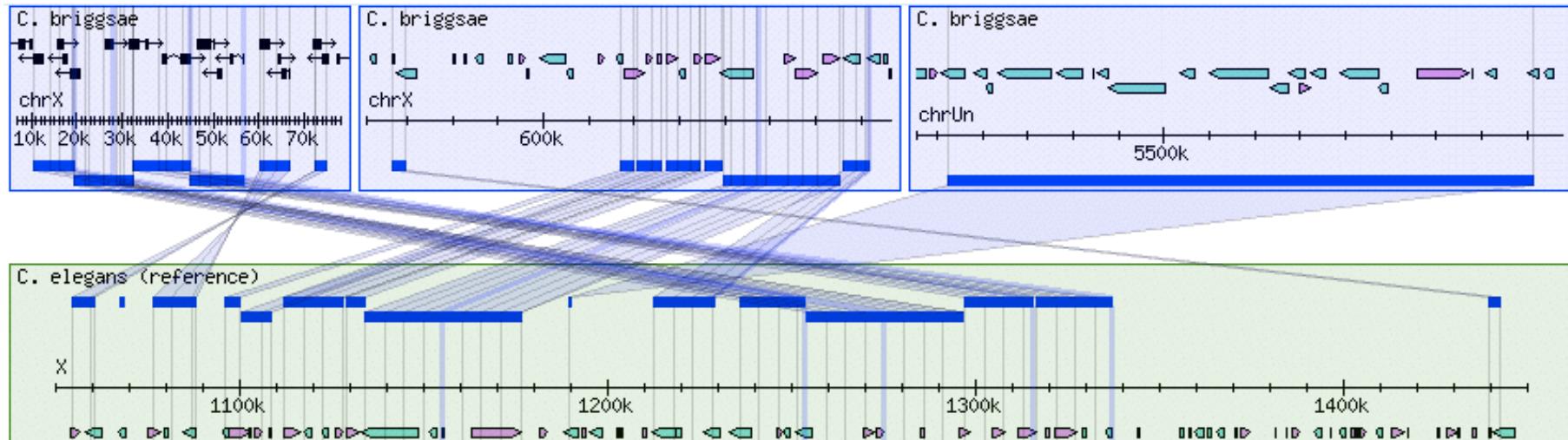
Générations successives



au sein d'une espèce: polymorphisme
entre espèces: divergence



2- Les génomes se réarrangent

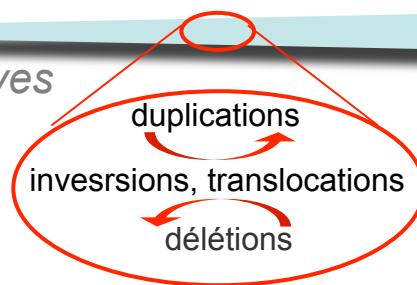


Thomas H. Morgan

Variations structurelles

Evolution des cartes génétiques

Générations successives



Perte de
synténie et
redondance



au sein d'une espèce: polymorphisme structurel, CNV
entre espèces: perte de synténie

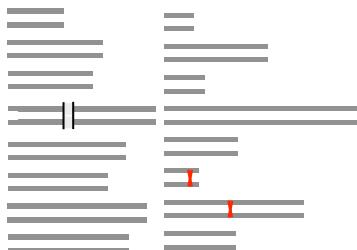
Point mutations and chromosomal alterations, two bases of genome evolution

2- Grandes altérations chromosomiques (changements d'ordre et/ou nombre de copies des gènes)



1576 inversions

total de **154 millions de nucléotides impliqués (~ 5 %)**



2 translocations

total de **680 000 nucléotides impliqués (~ 6 %)**

1- Mutations « ponctuelles » (changement de séquence)



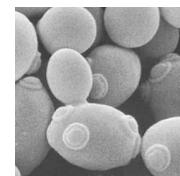
**35 millions de SNPs
5 millions indels**

~ **40 millions de nucléotides de divergence de séquence (~ 1.3 %)**



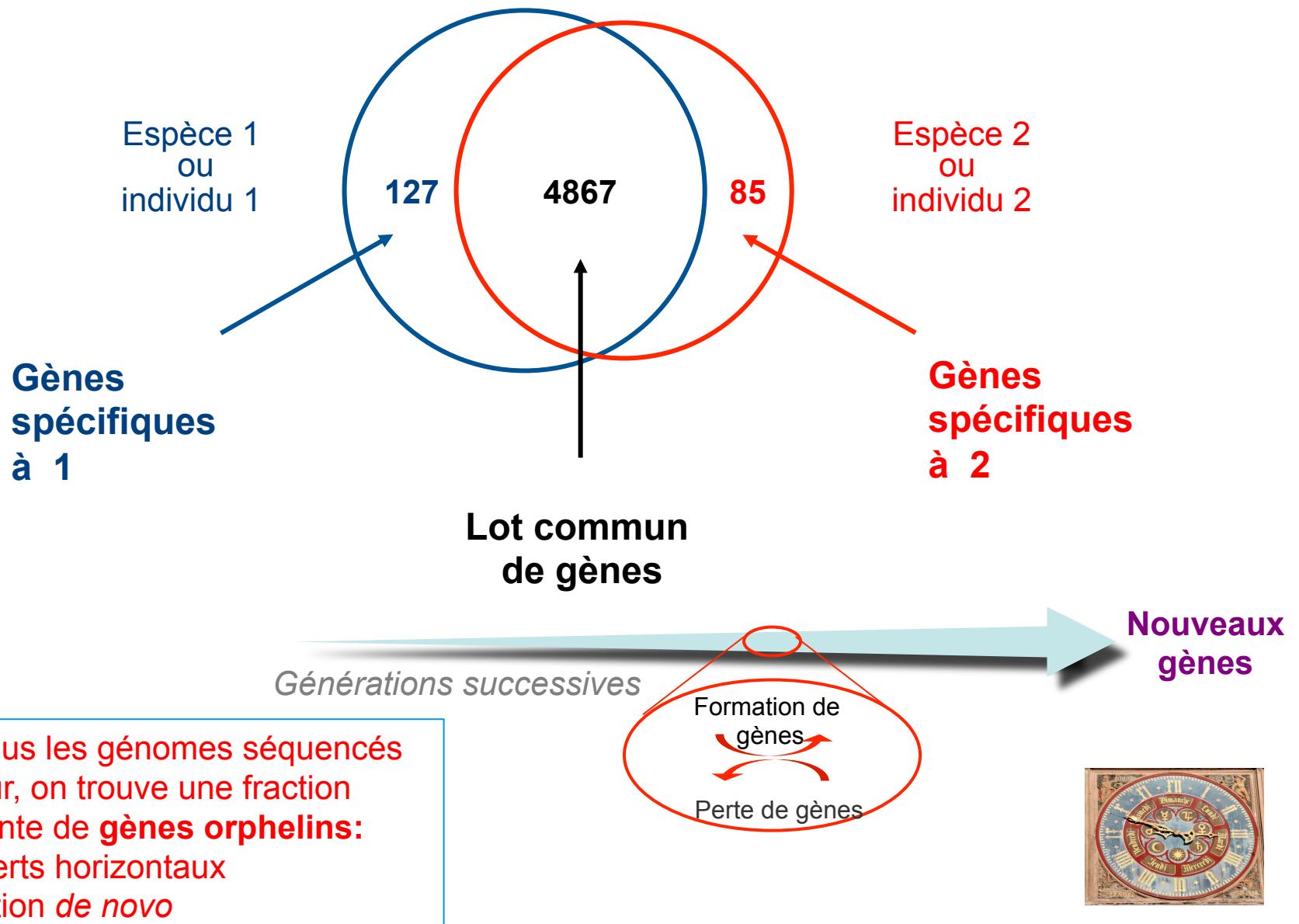
Saccharomyces cerevisiae

~ **4 millions de nucléotides de divergence de séquence (~ 33 %)**

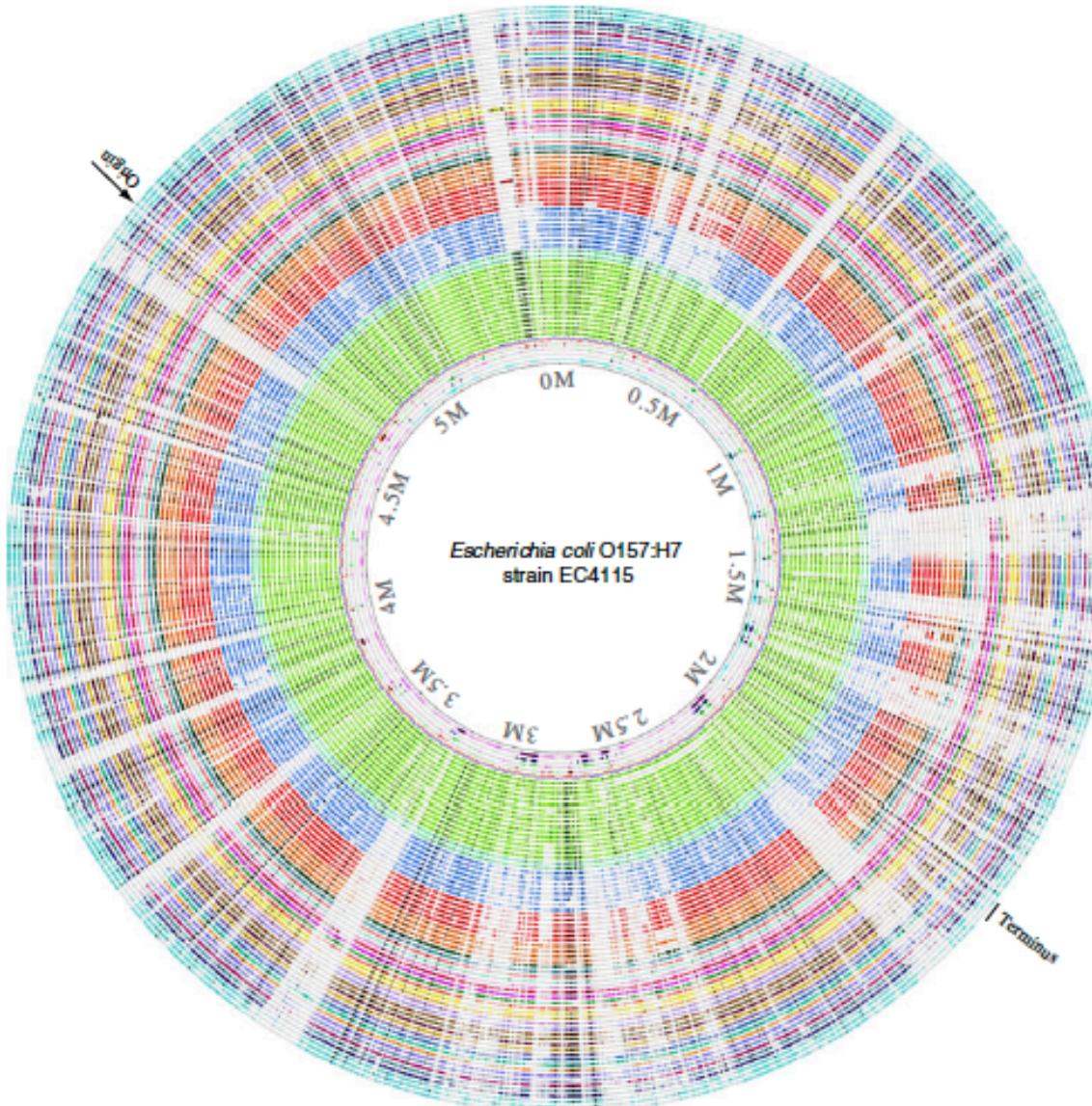


Saccharomyces uvarum

3- Des gènes apparaissent et disparaissent



Core-genome and pan-genome



core genome: essential cell functions

represents ~ **20 %** of each genome

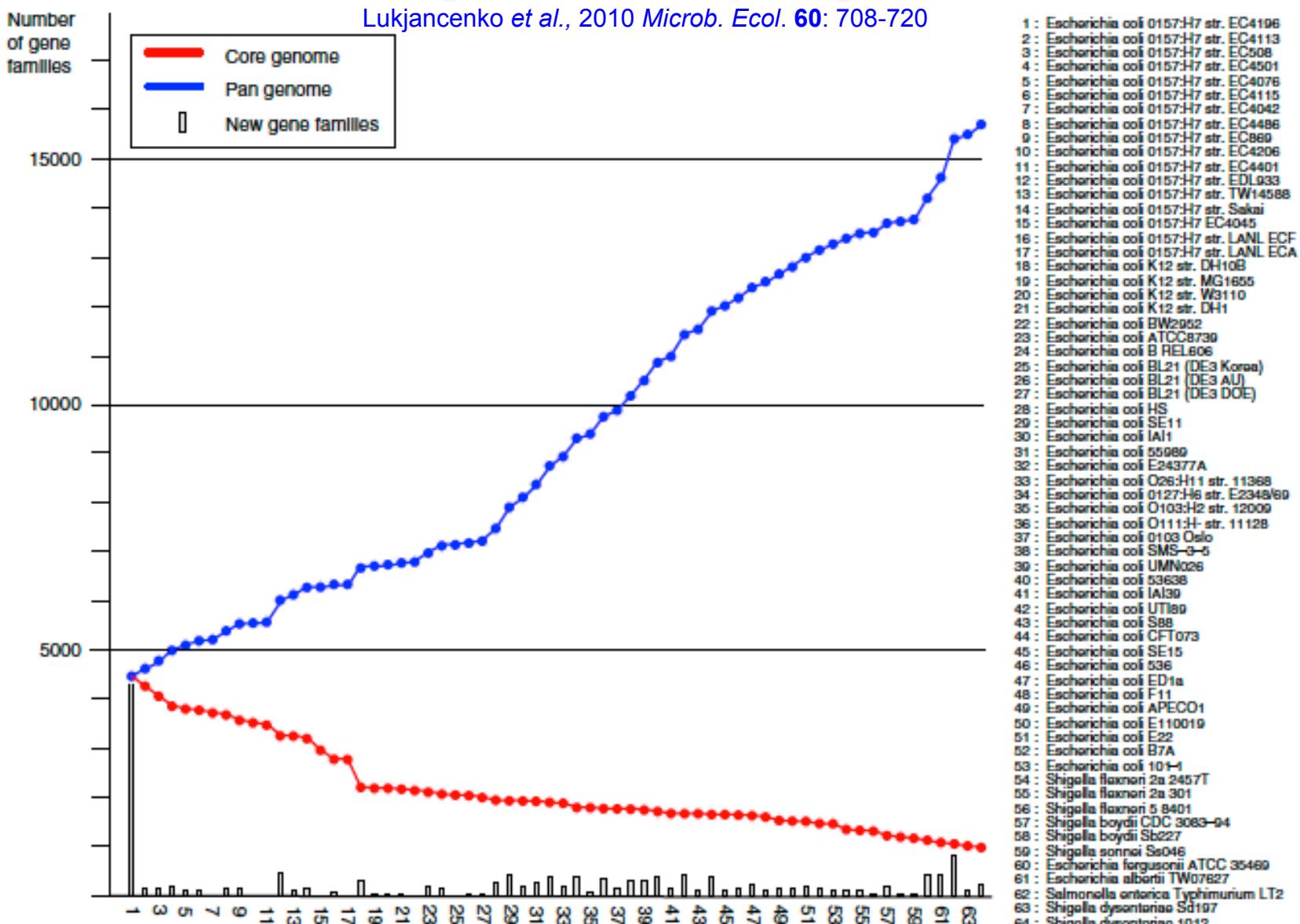
accessory genome: *fitness* and adaptation to different environments

Represents **80 %** of each genome

In *Bacteria*, accessory genes tend to co-localize on the genome, forming **genomic islands**.

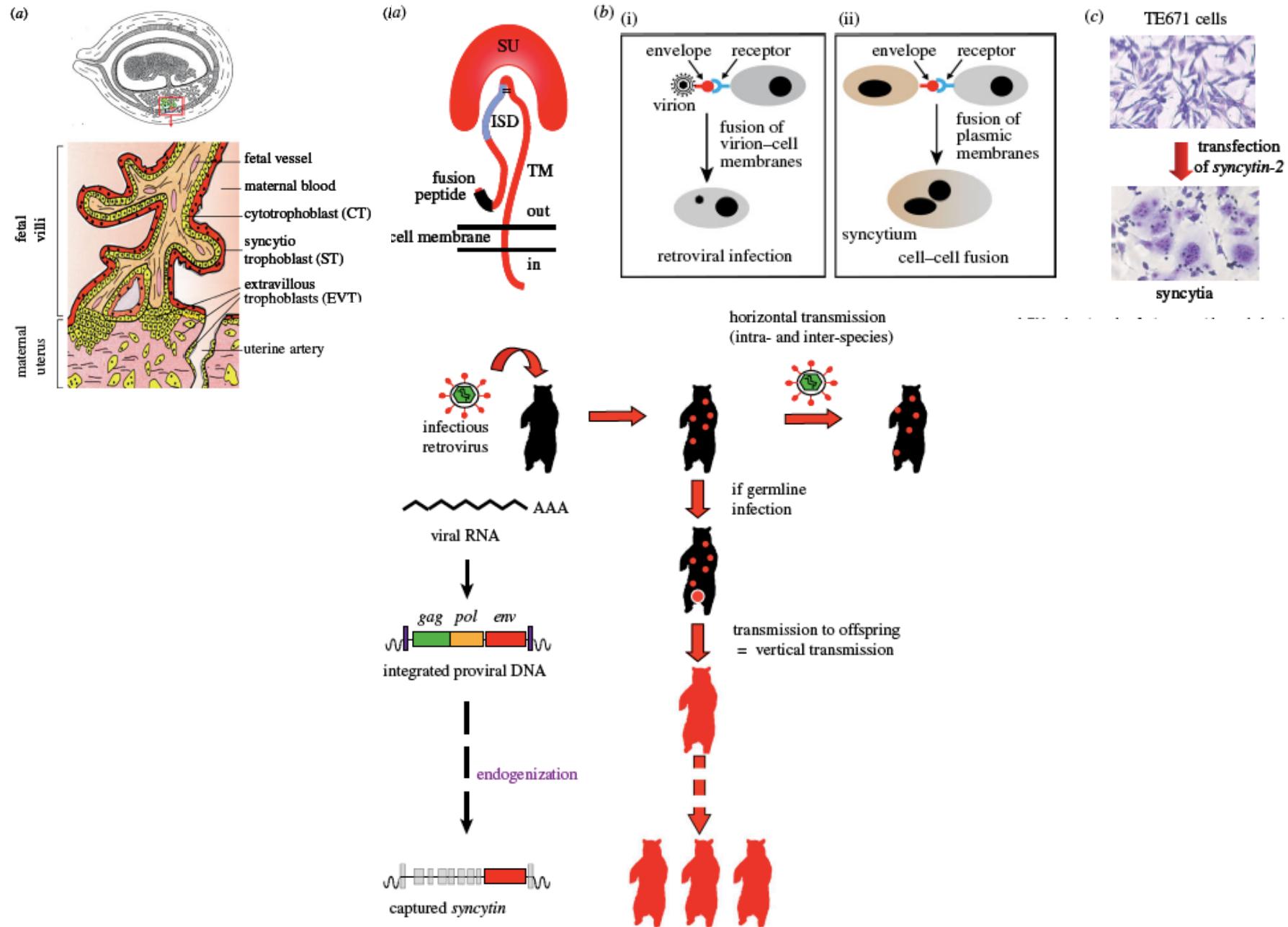
Core-genome and pan-genome

Lukjancenko et al., 2010 *Microb. Ecol.* 60: 708-720

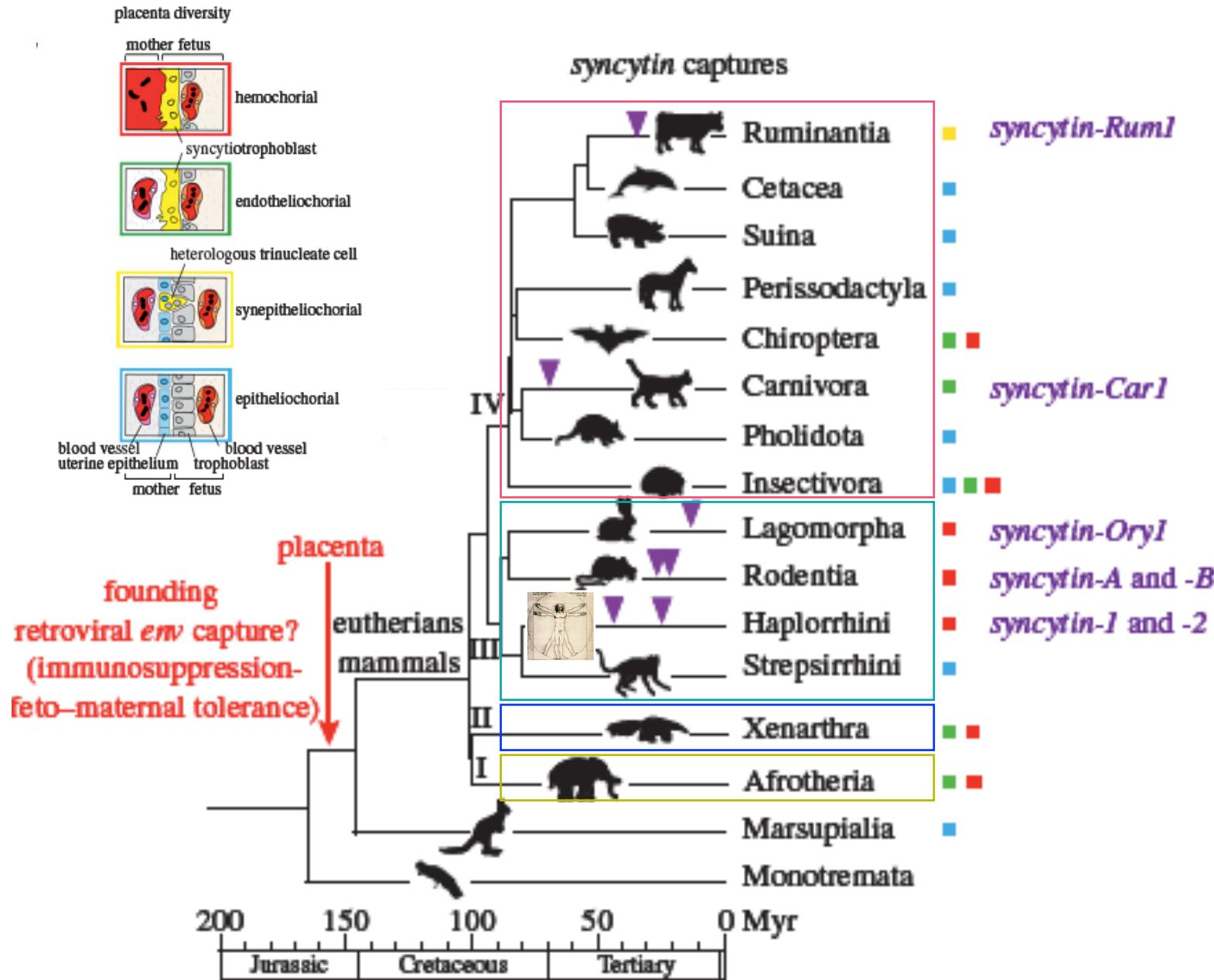


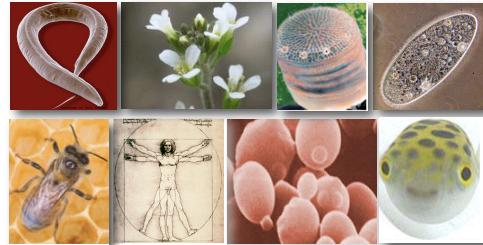
With 61 *E. coli* genomes sequenced, the pan-genome comprises 15 574 gene families of which only 993 (6 %) are present in all isolates (core-genome).

Important example of foreign gene acquisition: syncytines



A founding event followed by multiple similar events

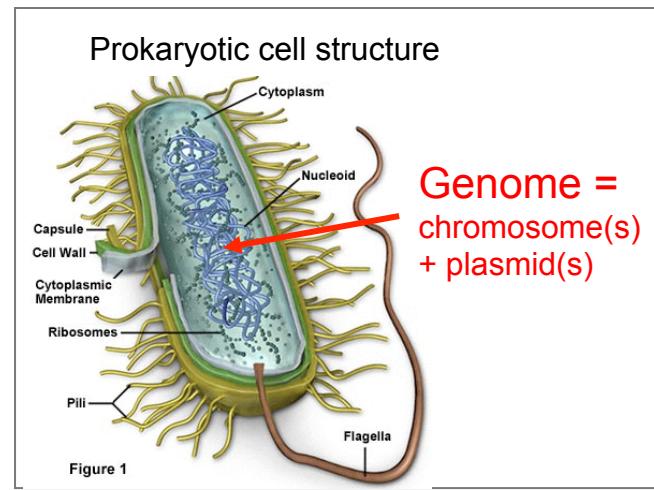




➤ Le monde des eucaryotes, vu par la génomique comparative

Les mécanismes moléculaires de l' évolution des génomes eucaryotes

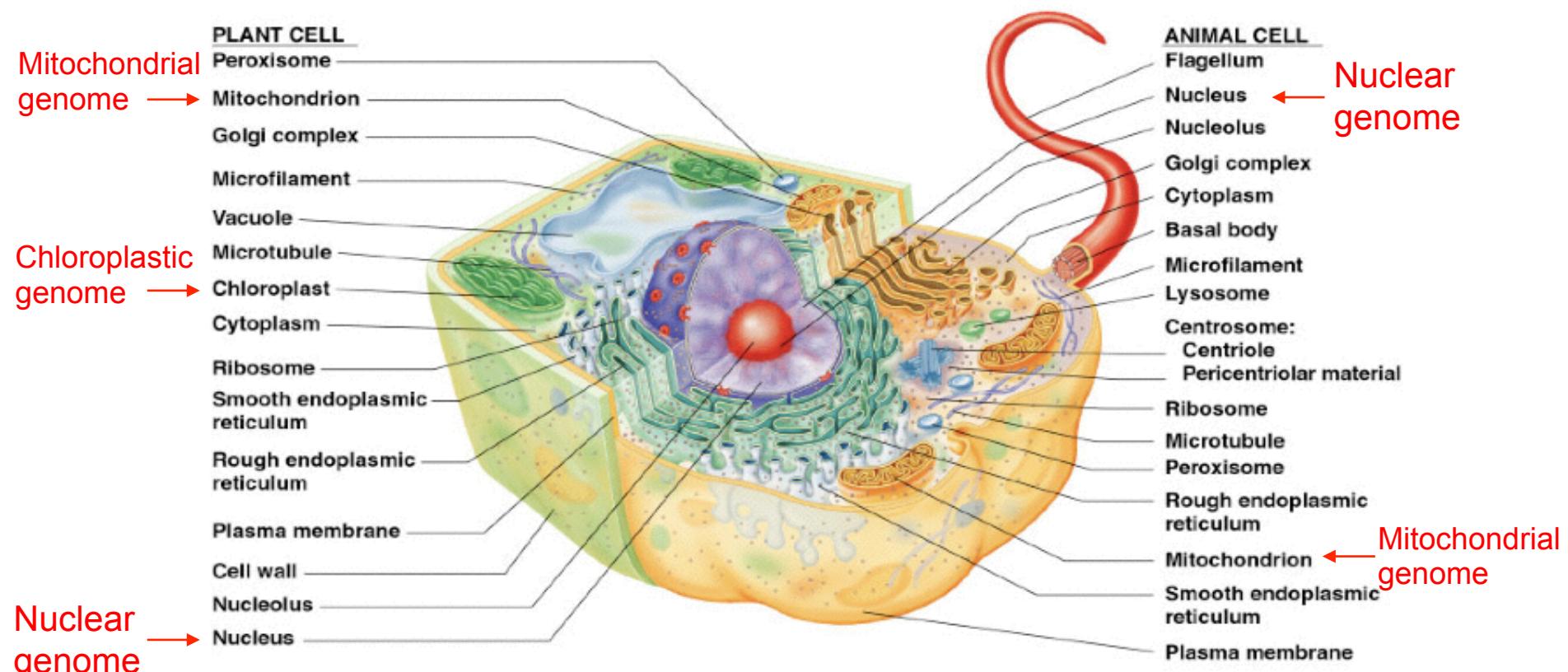
La formation *de novo* de gènes, *frameshifts* et ARNs



Eukaryotic cells

Genome =
nuclear chromosomes (+ plasmids)
+ mitochondrial genome (+ plasmids)
+ chloroplastic genome

Multiple membranes:
cell compartments

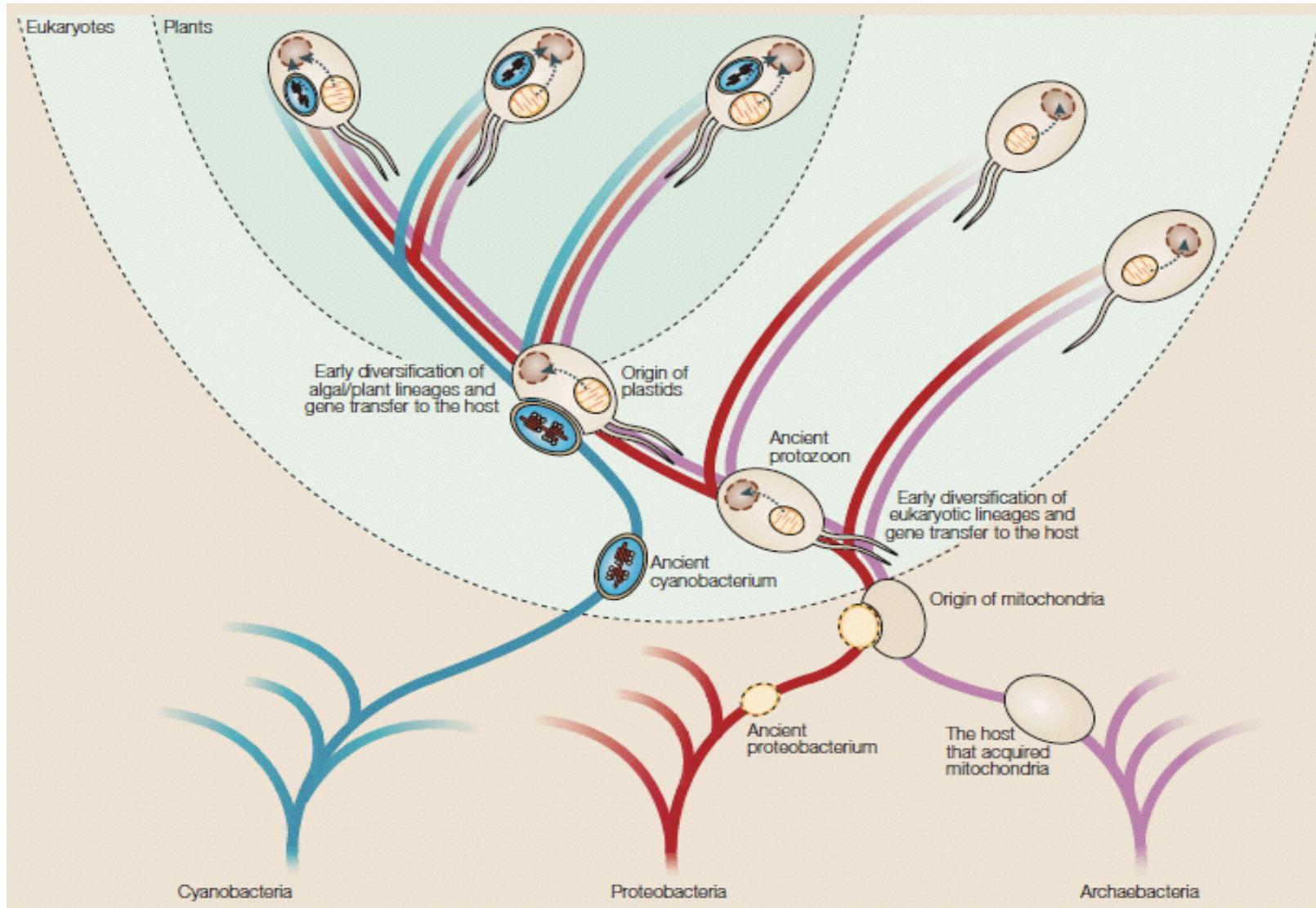


(a) Highly schematic diagram of a composite eukaryotic cell, half plant and half animal

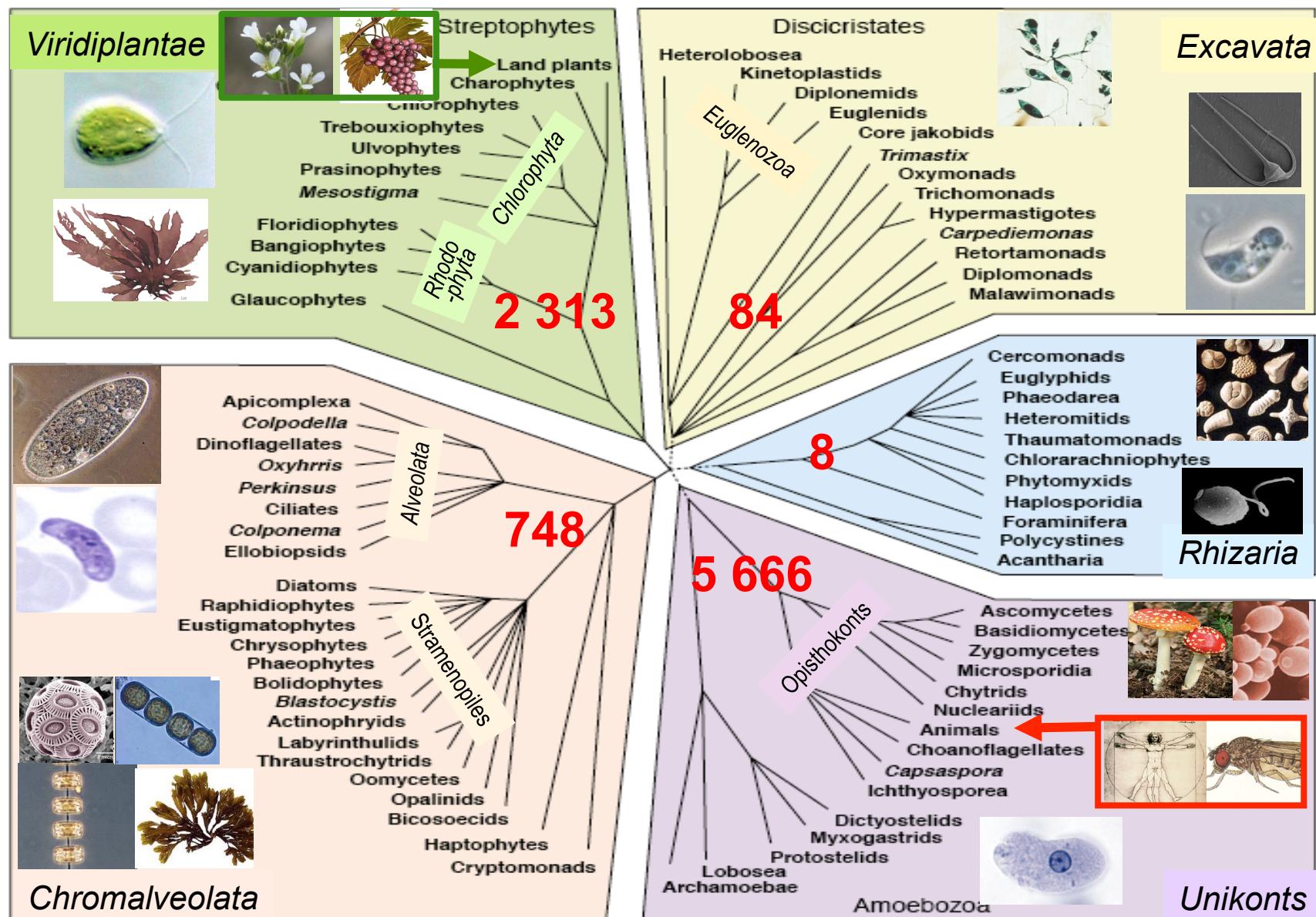
Copyright © 2004 Pearson Education, Inc., publishing as Benjamin Cummings.

Origin of eukaryotes: endosymbioses

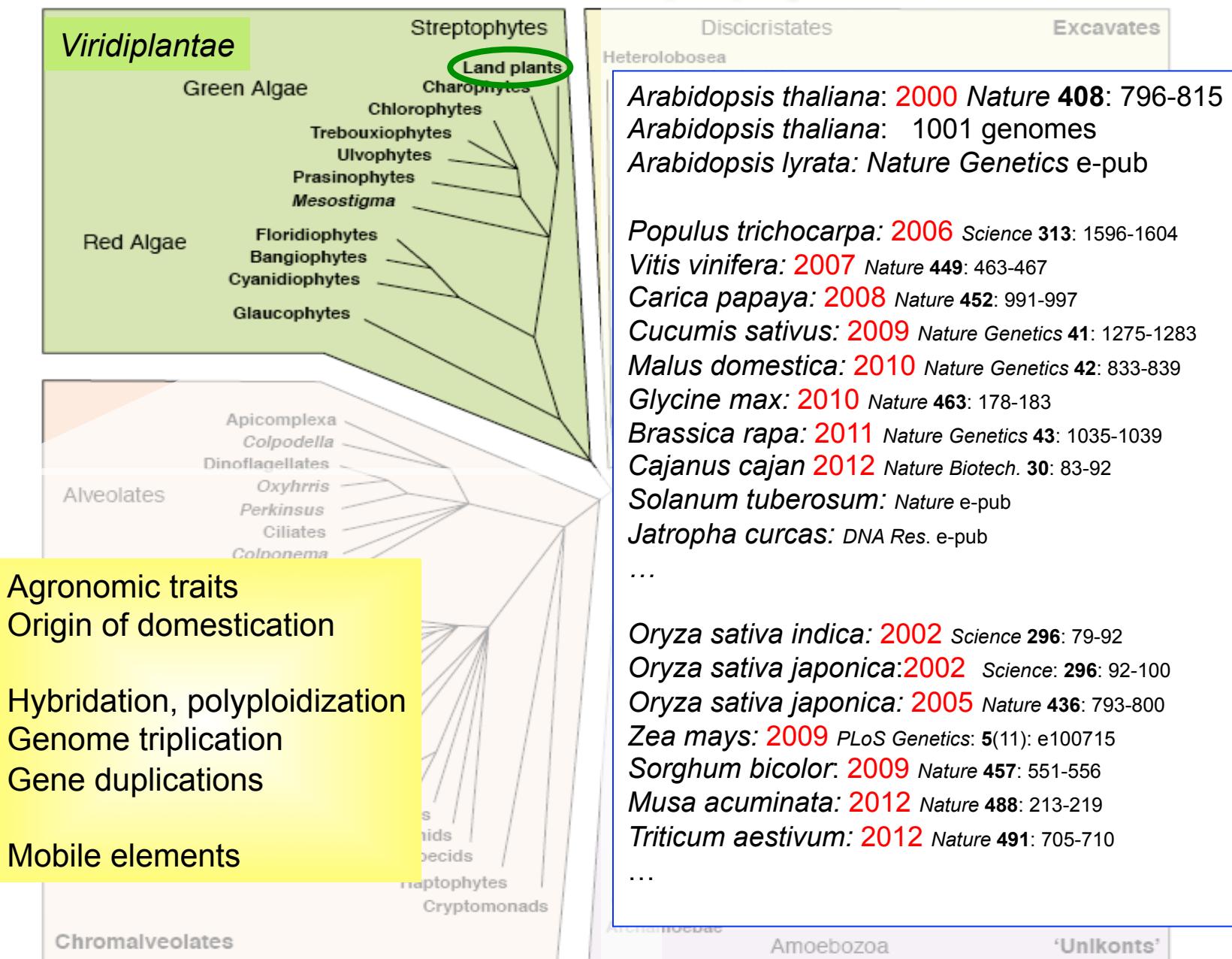
Timmis et al. (2004) *Nature Reviews Genetics* 5: 123-135



Genomics of the eukaryotes

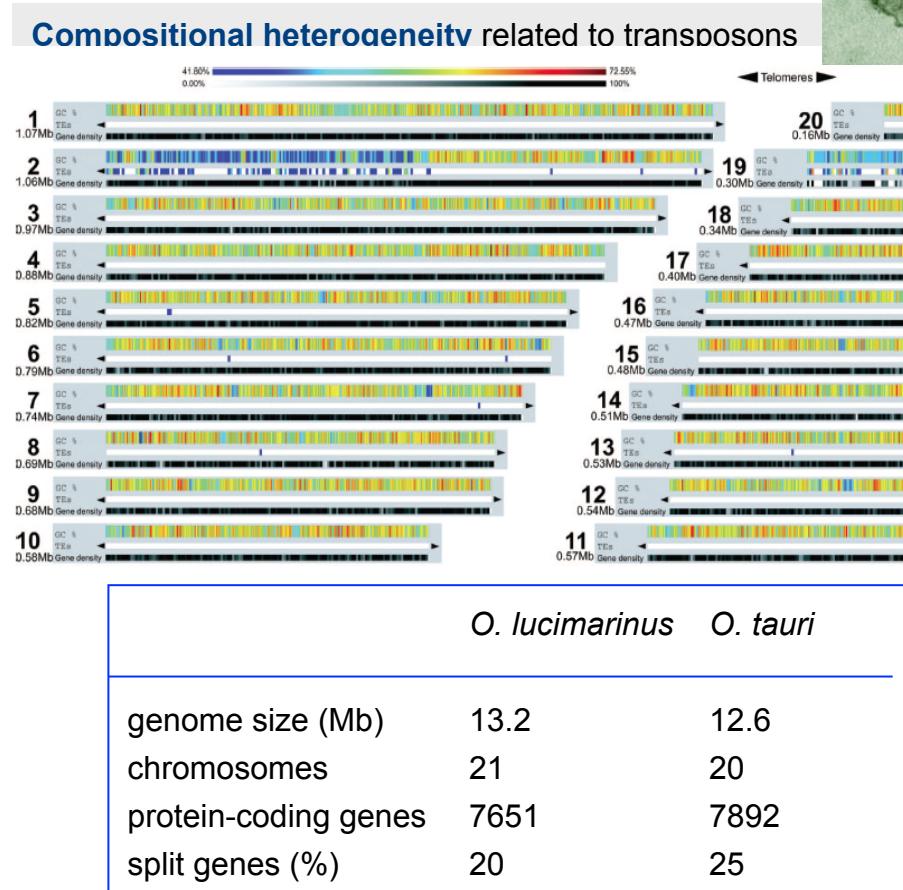


Genomes of Streptophyta

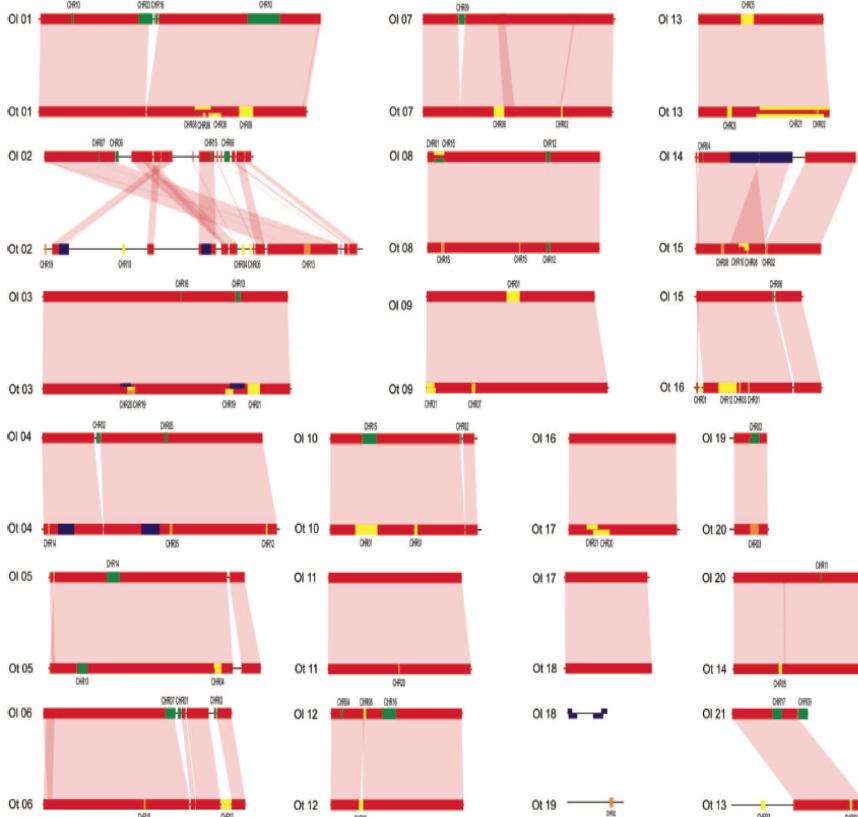


Genomes of unicellular eukaryotes: *Viridiplantae*

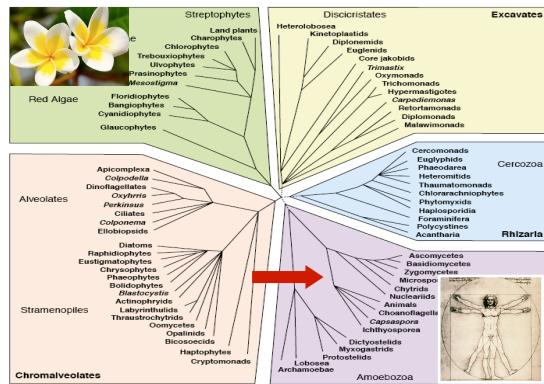
The genome of *Ostreococcus tauri*
Derelle et al., 2006 PNAS 103: 11647-11652
the smallest free living eukaryote



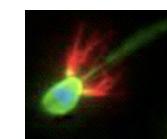
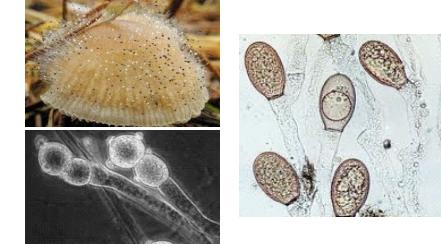
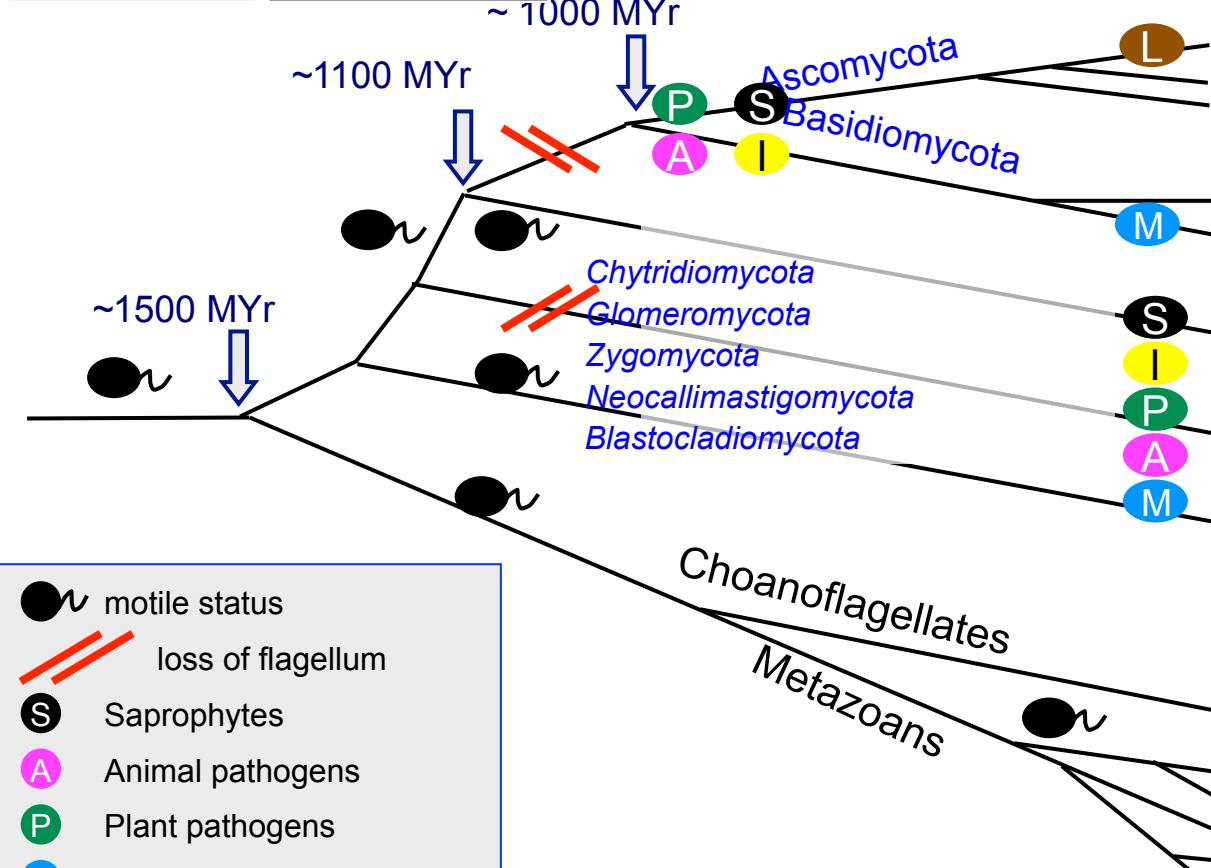
The genome of *Ostreococcus lucimarinus*
Palenik et al., 2007 PNAS 104: 7705-7710



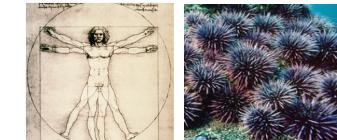
- Multiple mechanisms contribute to species divergence, act differently on different chromosomes.
- **Horizontal gene transfer** altering cell-surface characteristics.
- Numerous **gene fusions** 330 (O.t.), 348 (O. l.) of which 137 are common to both species.
- Numerous (20) genes for selenocysteine-containing proteins (TGA codons).



The world of Opisthokonts



Monosiga brevicollis



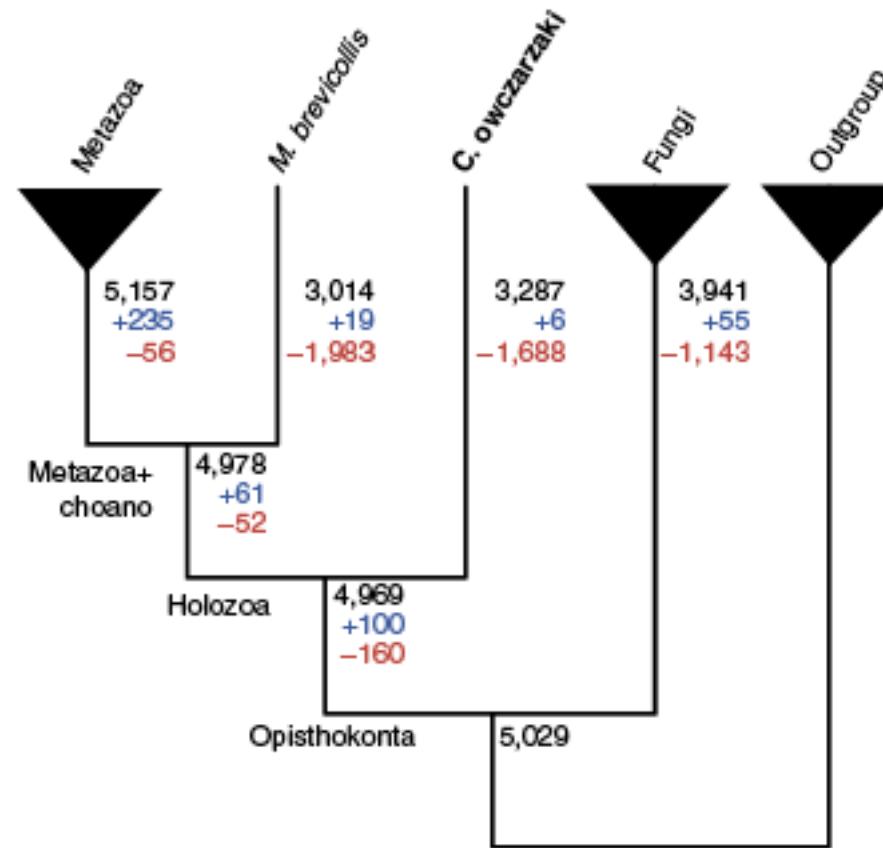
Deuterostomia



Ecdysozoa
Lophotrochozoa

The origin of metazoan protein domains

Total number of Pfam domains and numbers of domain gain (+) and loss (-) events inferred from Dollo parsimony



Gain and loss of protein domains within the Opisthokonta

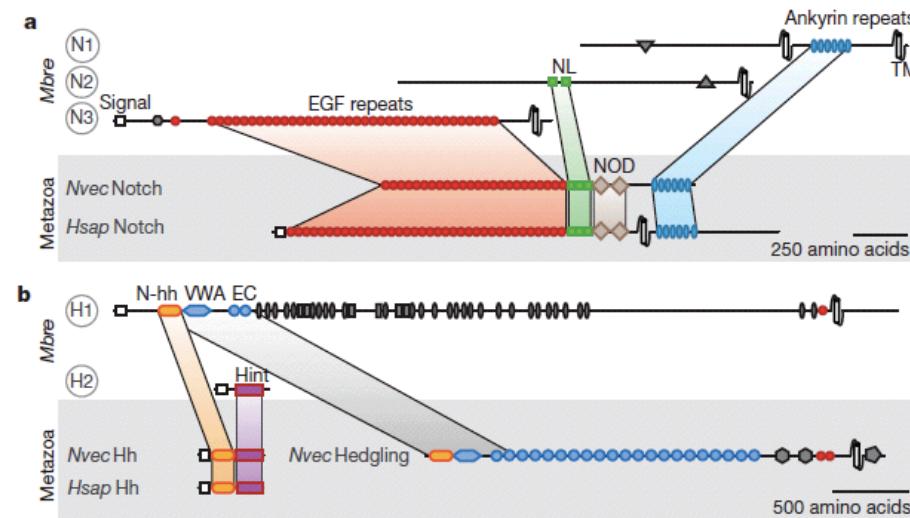
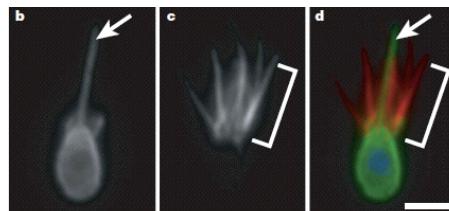
The origin of metazoan protein domains

Monosiga brevicollis
(Choanoflagellate)

King et al., 2008, *Nature* 451: 783-788

41.6 Mb

9 171 protein-coding genes



Examples of protein domain fusions

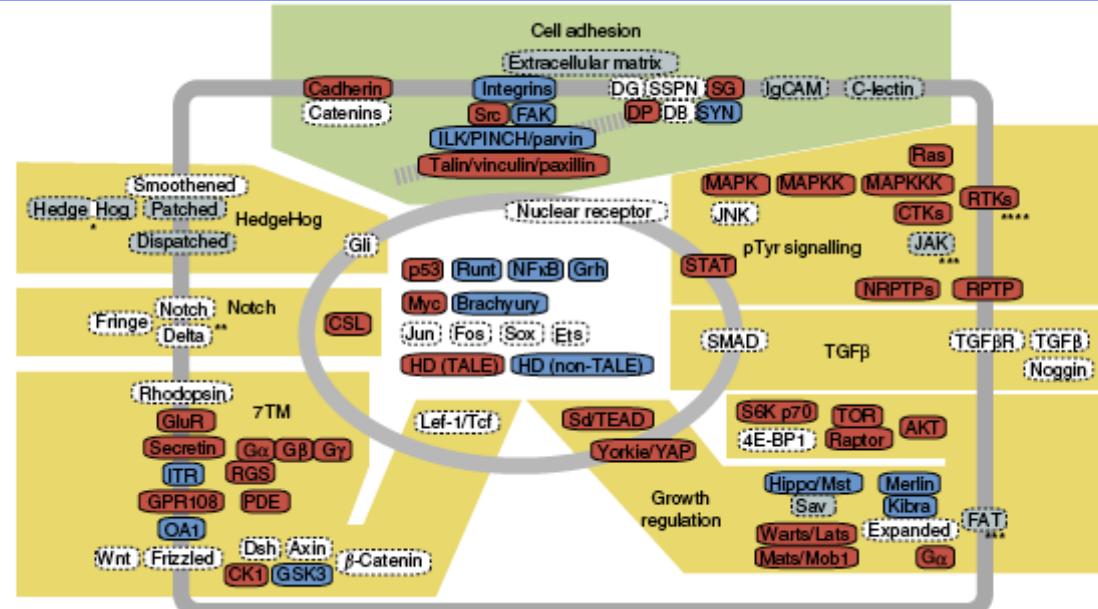
Capsaspora owczarzaki
(Filasterean)

Suga et al., 2013, *Nature Commun.* 4: 2325

28.0 Mb

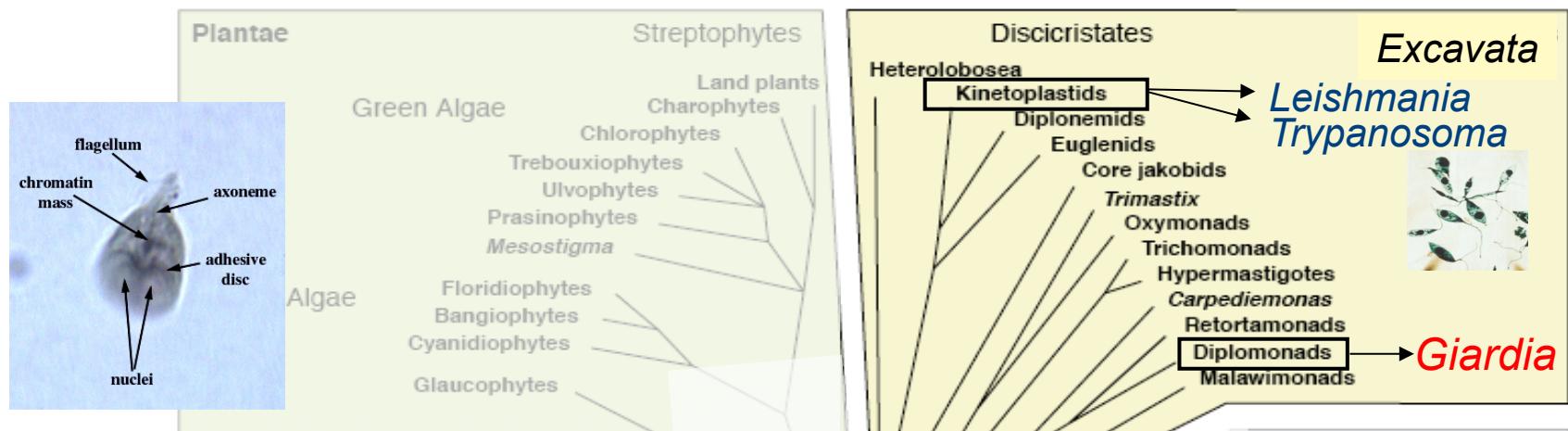
8 657 protein-coding genes

- Present in *M. bre.* and *C. owc.*
- Present in *C. owc.* only.
- Present in *M. bre.* only.
- Absent in *M. bre.* and *C. owc.*



Traces of protein components of major metazoan cell adhesion complexes (green) and various signalling pathways (yellow).

Genomes of unicellular eukaryotes: Excavata

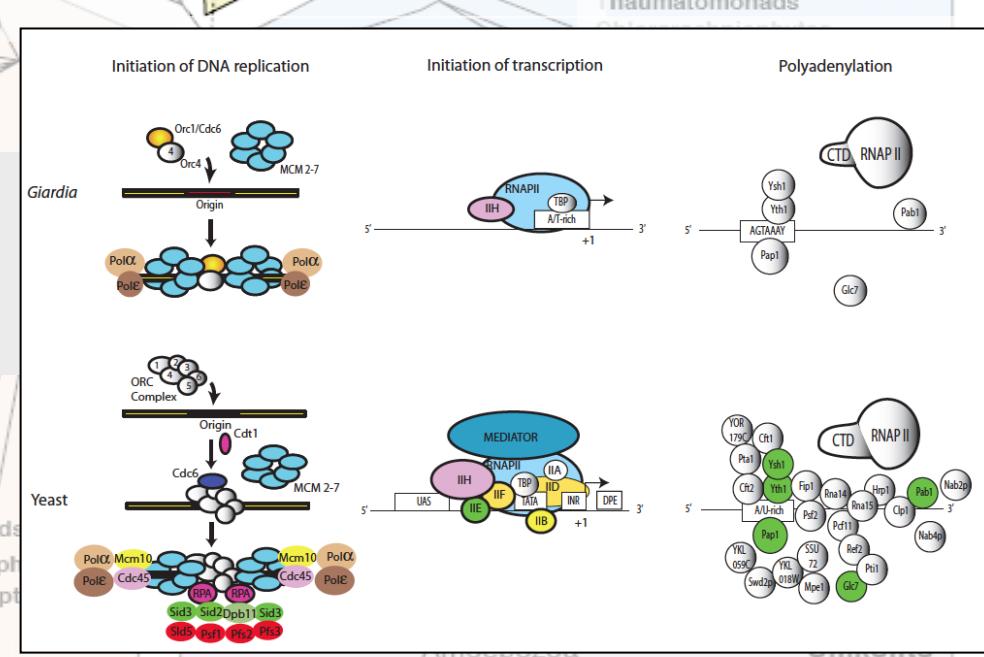


The genome of *Giardia lamblia*

Morrison et al., 2007 Science 317:1921-1926

human intestinal parasite, flagellated trophozoites attach to epithelial cells
two diploid nuclei, no mitochondria, no peroxisomes

- Genome ~11.7 Mb, 5 chromosomes, 6470 annotated CDS, very few introns (4)
- Low degree of heterozygosity (0.01% between the 4 genomes)

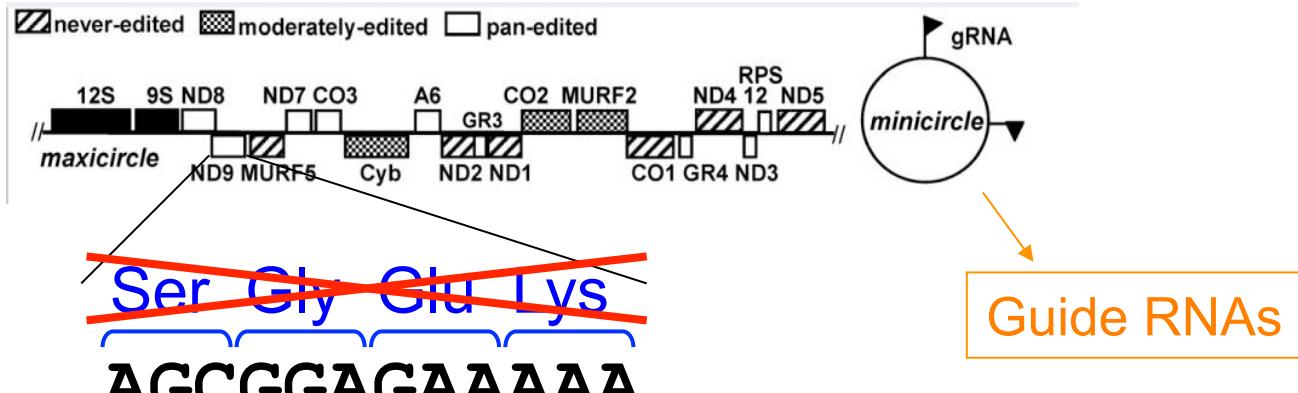


➤ Evolution is often regressive

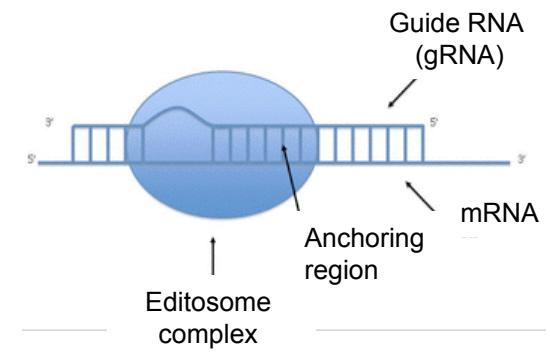


1983 Kinetoplasts of *Trypanosoma*

RNA editing

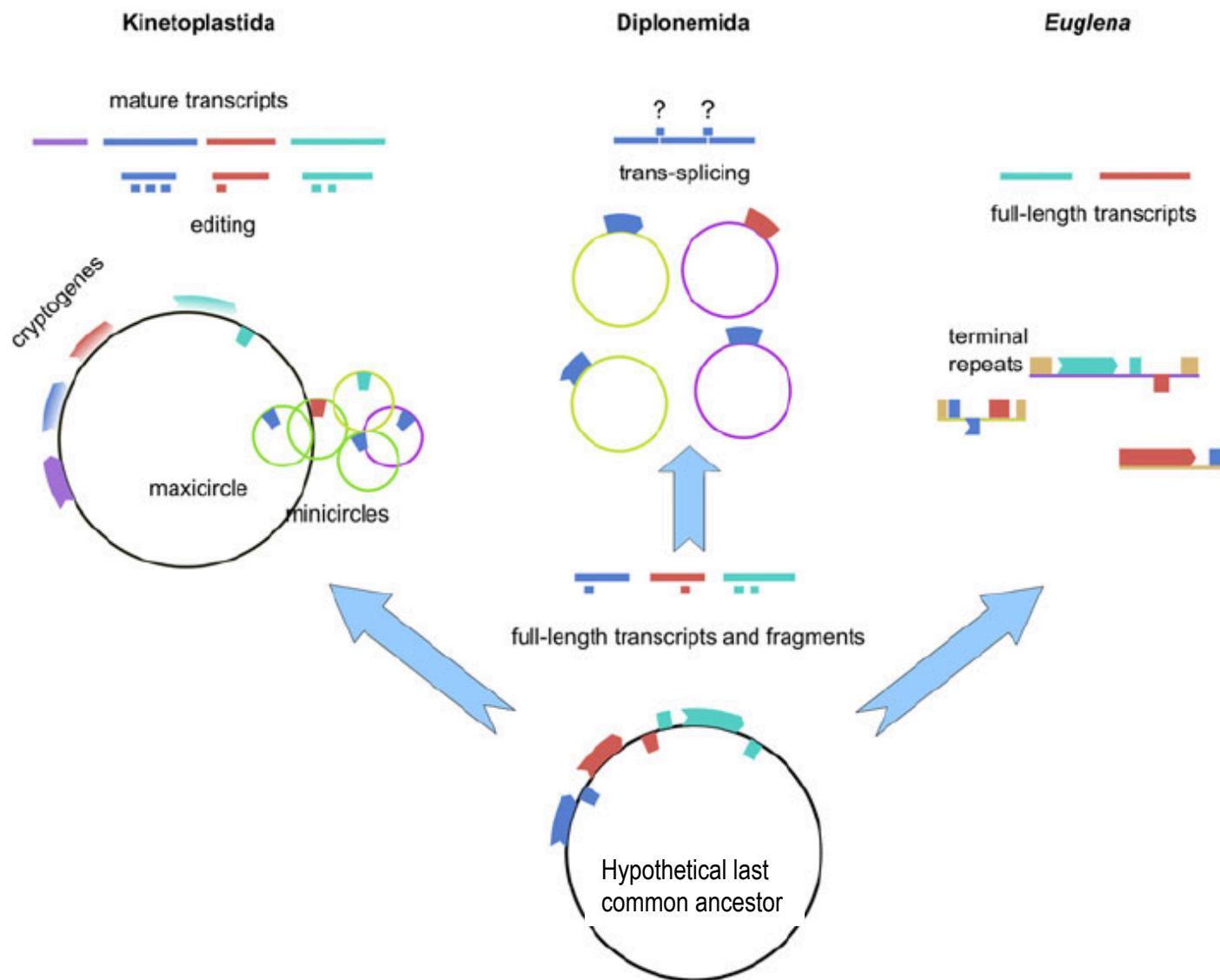


AuGuuuCGuuGuAGAuuuuuuAuuAuuuuuuuuAuuA
Met Phe Arg Cys Arg Phe Leu Leu Phe Phe Leu Leu

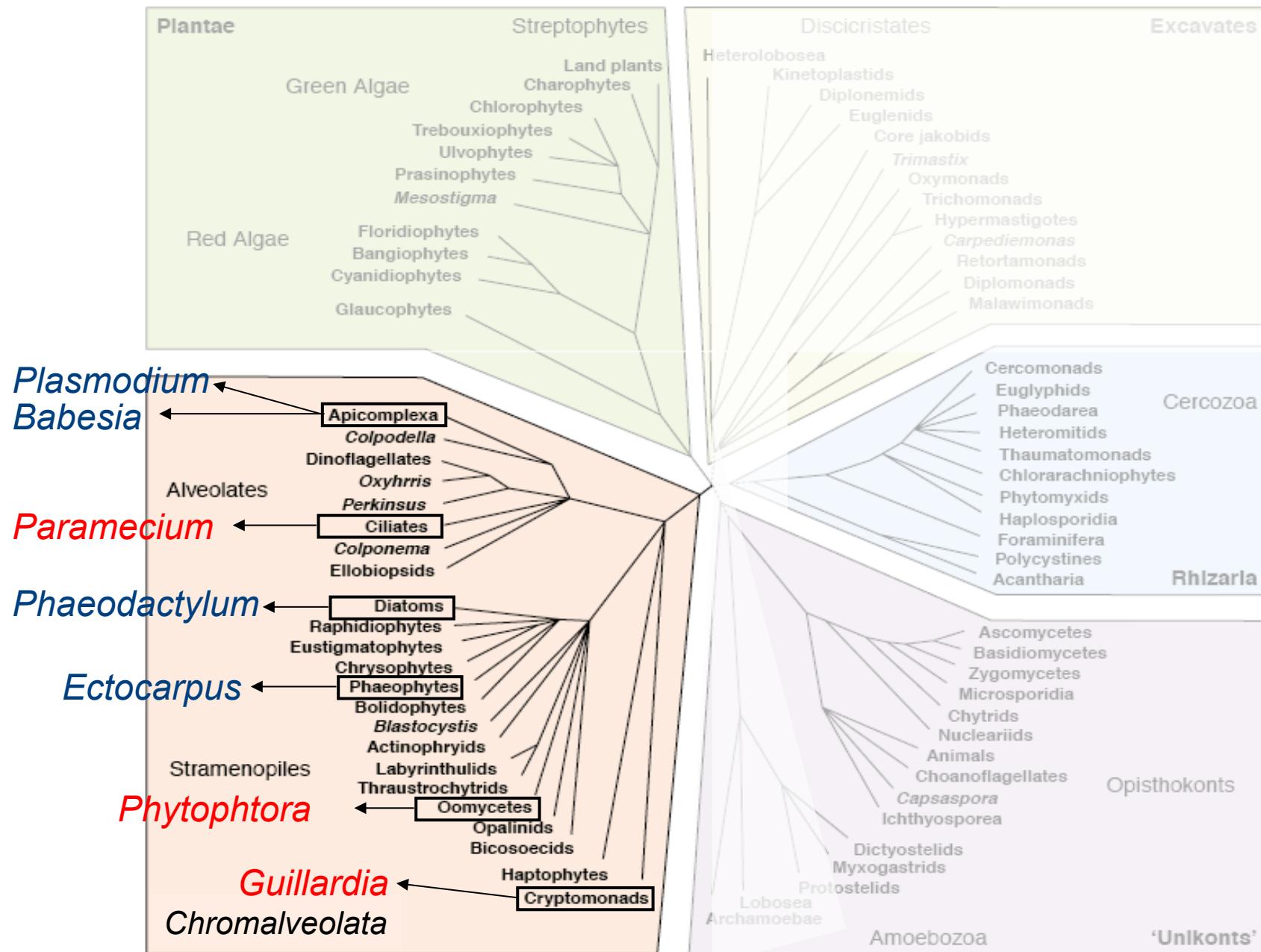


Germline mitochondrial genomes in *Euglenozoa*

Flegontov et al., (2011) Curr. Genet. 57: 225-232



Genomes of unicellular eukaryotes: Chromalveolata



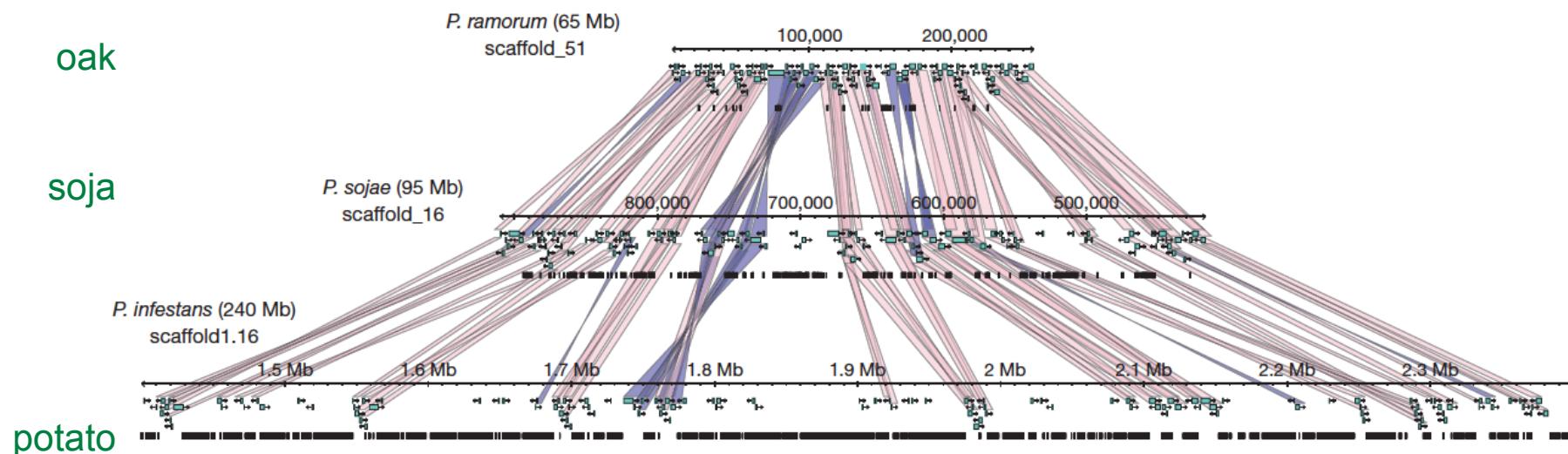
The genomes of *Phytophthora infestans* (*sojae* and *ramorum*)

Haas et al., 2009 *Nature* 461: 393-398

	<i>P. infestans</i>	<i>P. sojae</i>	<i>P. ramorum</i>
Genome size (Mb)	240	95	65
Scaffolds	4921	1810	2576
Repeat (%)	74	39	28
Protein-coding genes	17797	16988	14451



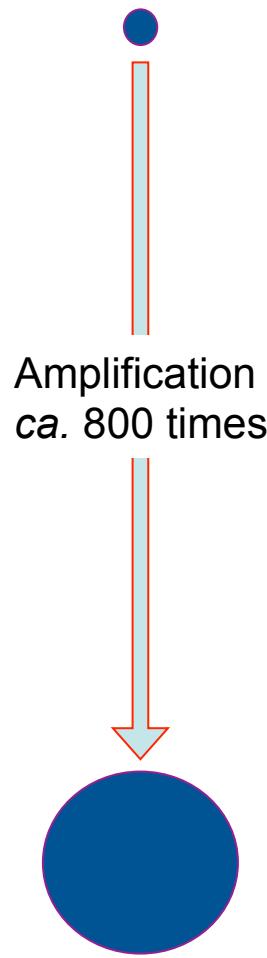
- **Conserved syntenic blocks** containing most common genes and few repeated DNA are separated by regions of repeated DNA with low gene density and no conservation of gene order --> **dynamic genomes**
- **Rapidly evolving effector genes** in non-conserved regions (modular secreted proteins, major types: RXLR and Crinkler, targetted to plant cells and responsible for necrosis, mostly species-specific)
- Effector **gene family expansion** in *P. infestans* associated with numerous mobile elements (helitron)



Genomes of unicellular eukaryotes: *Chromalveolata*

The macronuclear genome of *Paramecium tetraurelia*

Aury et al., (2006) *Nature* 444: 171-178 (697 scaffolds, totalling 72 Mb)



Micronucleus (2n) genome size
ca. 100 Mb
> 50 chromosomes

- Precise elimination of > 10 000 short, unique copy elements ---> **Reconstruction of functional genes**

RNA mediated process

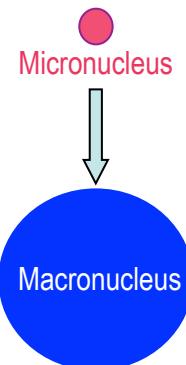
- Imprecise elimination of transposable elements and other repeated sequences
---> **Chromosome fragmentation, de novo telomere addition, internal deletions**

Macronucleus (polyploid)
genome size ca. 75 Mb
39 642 genes

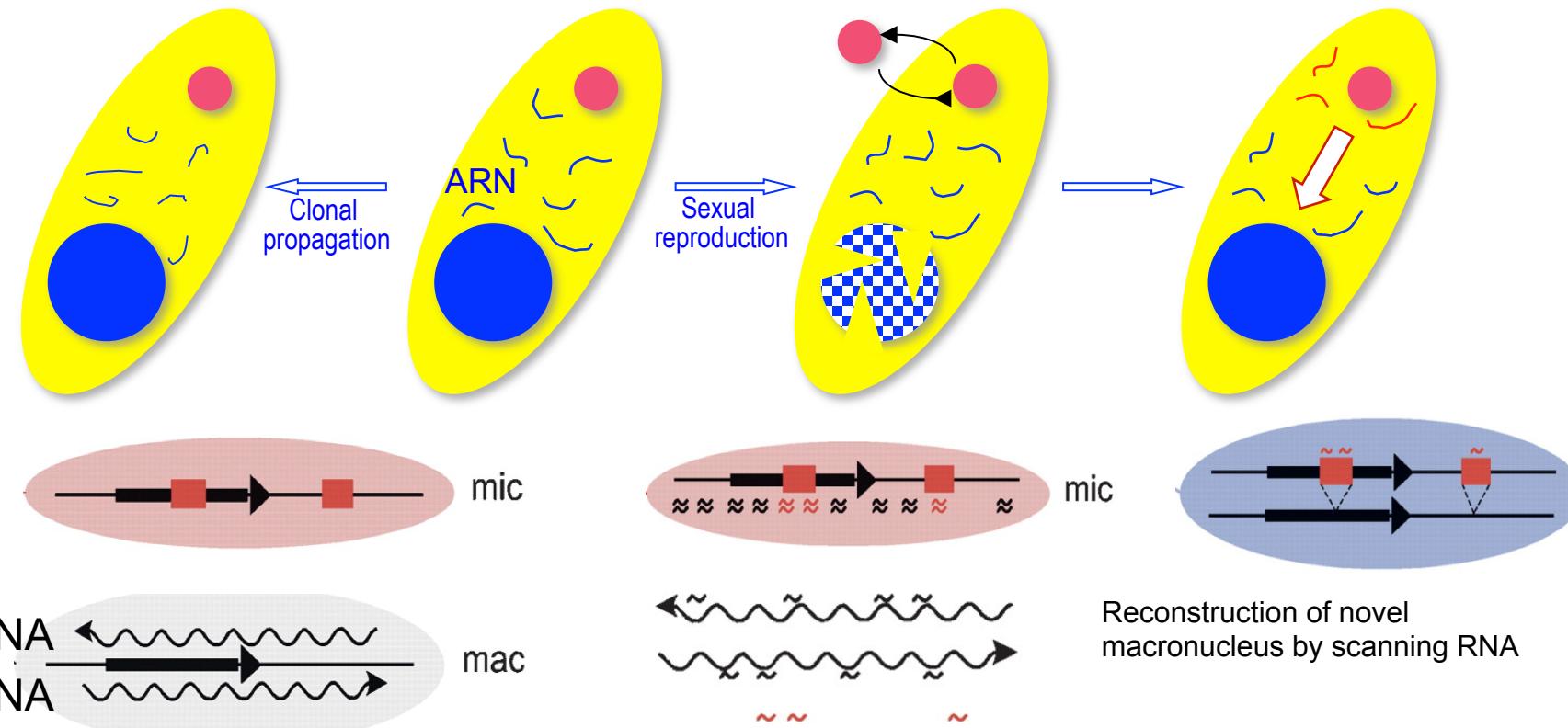
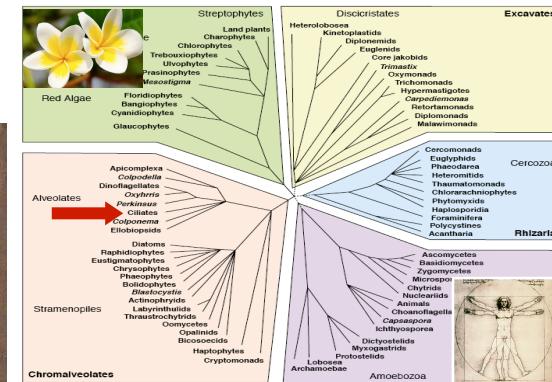
Reconstruction of a functional genome from the germline genome

Aury et al., (2006); Lepèvre G et al. (2008)

- Precise elimination of > 10 000 short, unique copy elements ---> **Reconstruction of functional genes RNA mediated process**
- Imprecise elimination of transposable elements and other repeated sequences ---> **Chromosome fragmentation, de novo telomere addition, internal deletions**



Paramecium tetraurelia

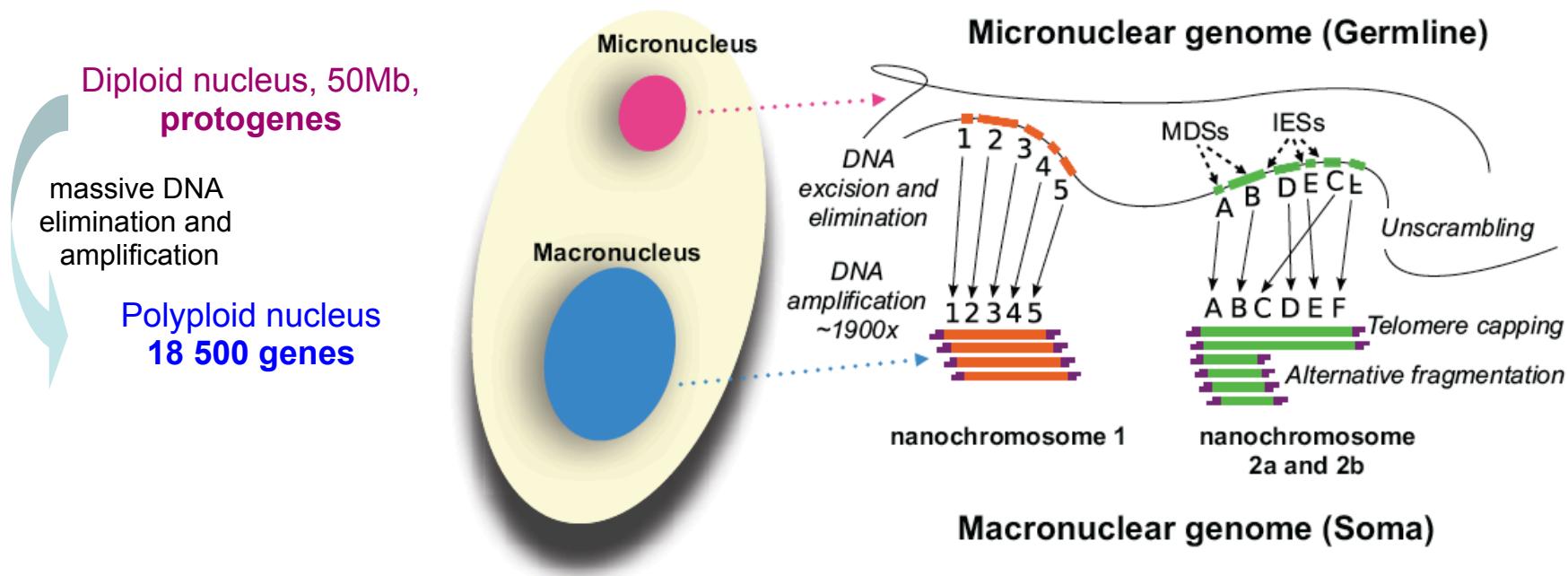




Reconstruction of a functional genome from the germline genome

Oxytricha trifallax

Stuart et al., 2013 PLoS Biology 11:1 e1001473



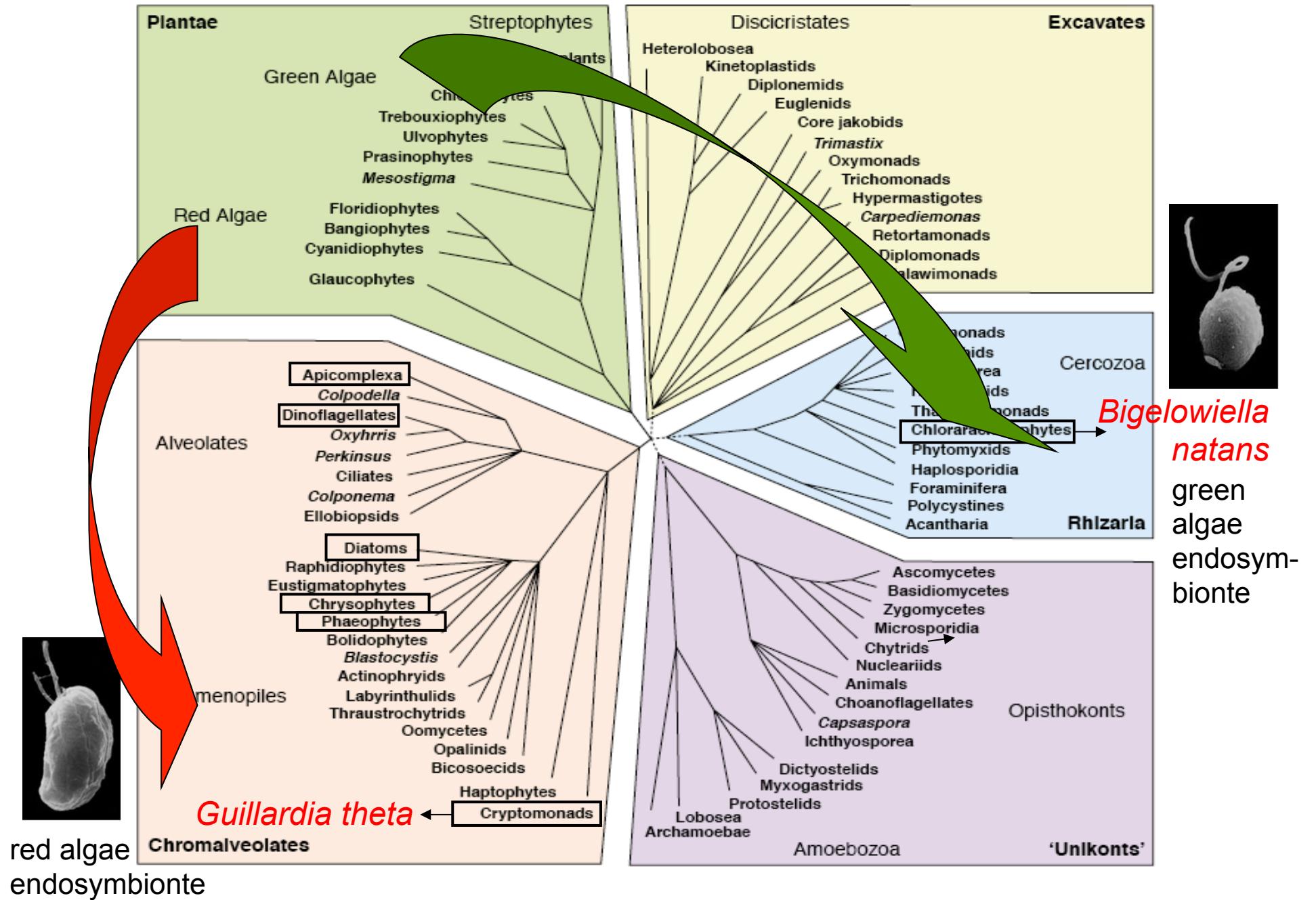
> 16 000 nanochromosomes (length 0,4 à 66 kb, mean 3,2 kb)

90 % correspond to a single gene (others to 2-8 genes)

Multiple isoforms, variable amplification (mean 2 000 copies) > polymorphism between individuals

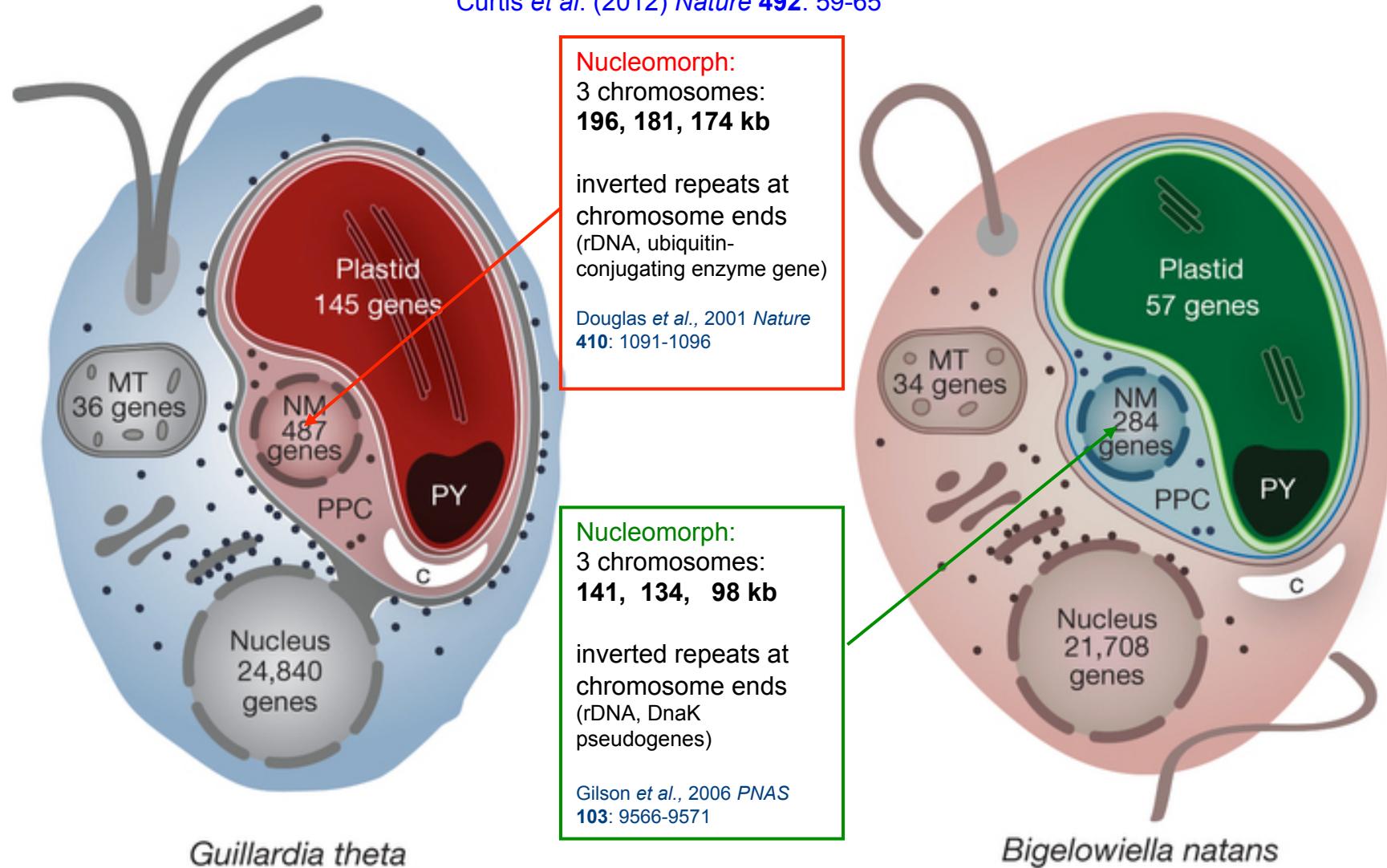
de novo generation of genes by reassociation from protogene elements. millions of telomeres

Four membrane plastids and nucleomorphs (double endosymbiosis)

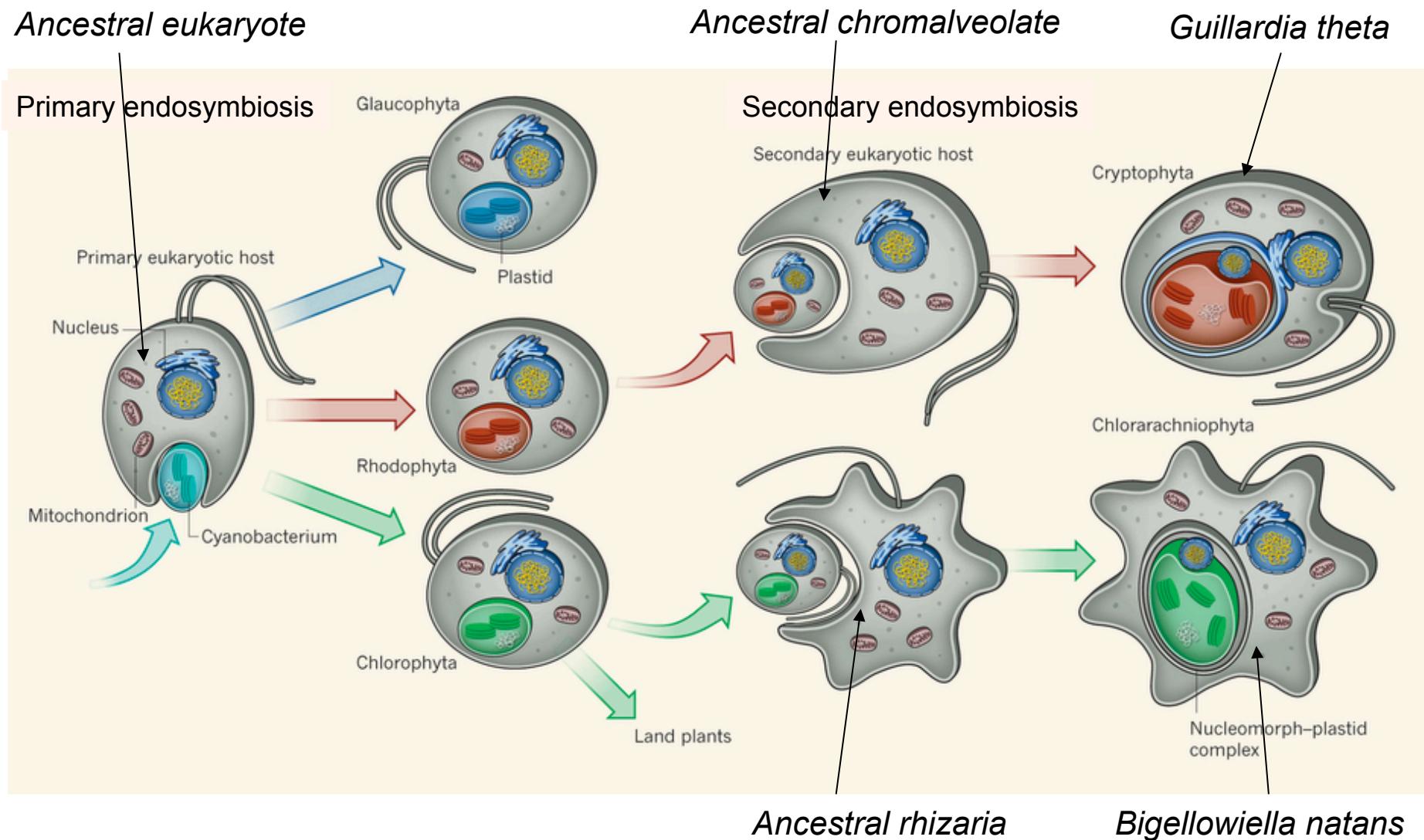


Algal genomes reveal evolutionary mosaicism and the fate of nucleomorphs

Curtis et al. (2012) *Nature* 492: 59-65



	<i>Guillardia theta</i>	<i>Bigelowia natans</i>
Genome size (Mb)	87.2	94.7
Split CDS (%)	80	86
Mean exons /gene	6.4	8.8
Mean intron size (nuc.)	110	184



The world of eukaryotes from genomics

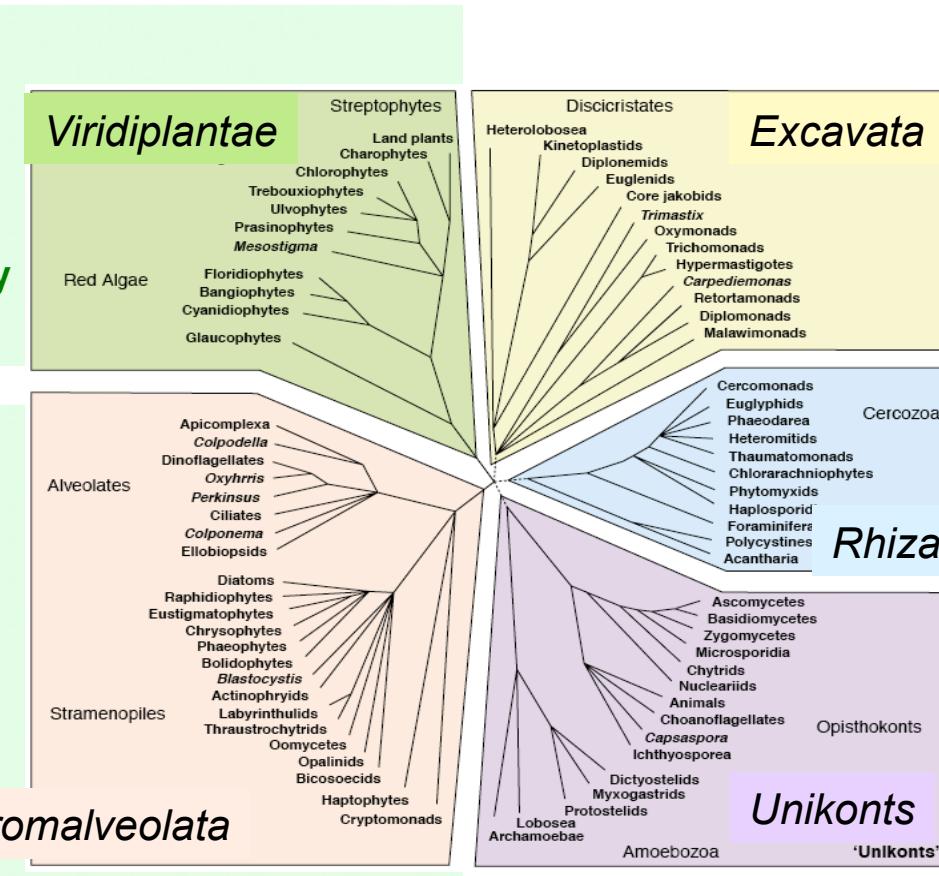
From small, compact genomes to massively duplicated genomes

Primary endosymbiosis (chloroplasts) followed by chloroplast loss

Primary endosymbiosis followed by chloroplast loss

Genome reconstructions (with RNA intermediate)

Recent endosymbiosys (red algae)



Massive RNA edition
Regressive evolution

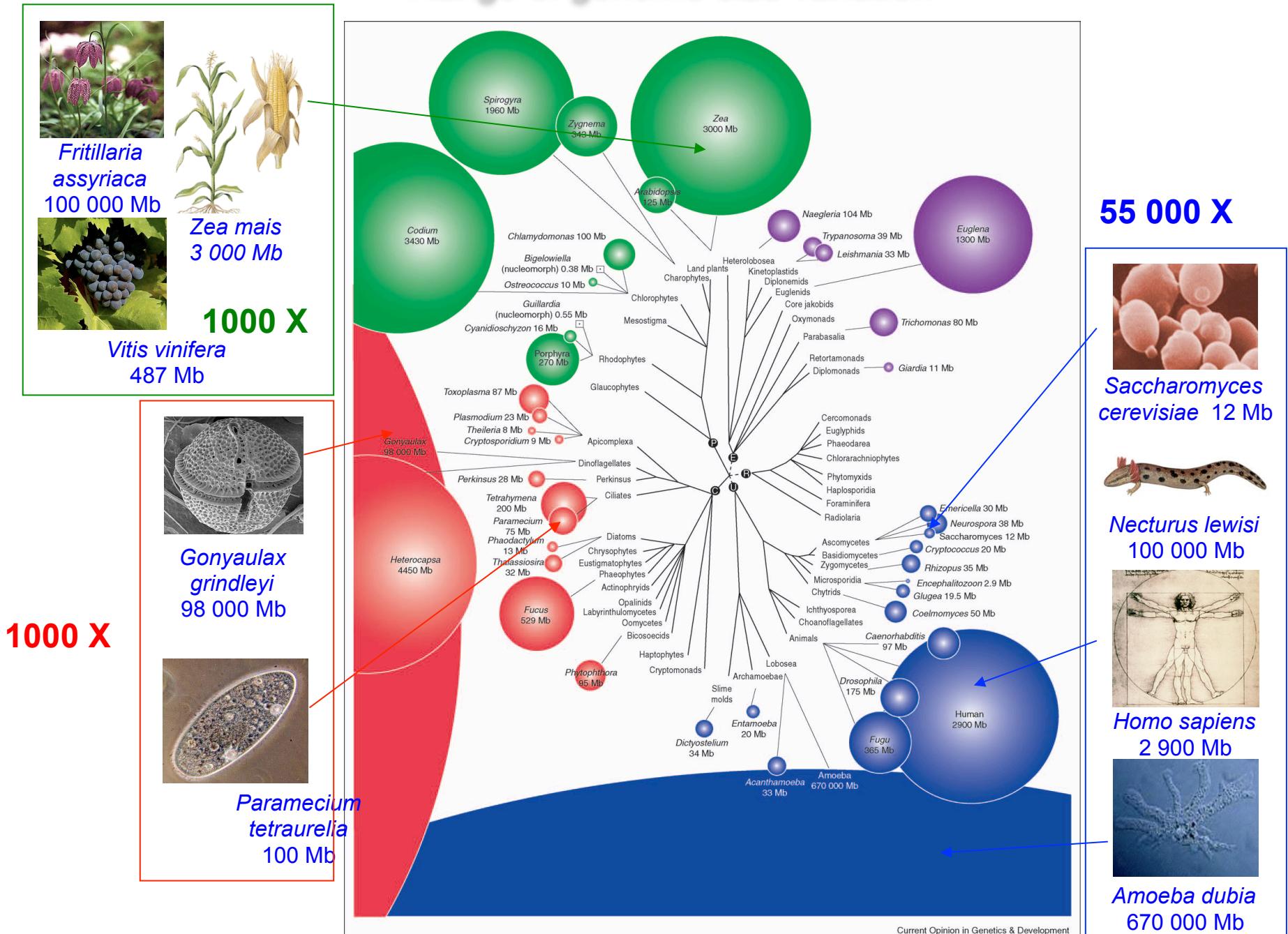
Recent endosymbiosis (green algae)

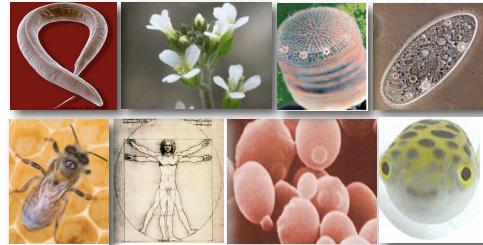
What else ?

From small, compact genomes to extremely large genomes

Gene duplication, genome duplication, intron gain and loss, horizontal gene acquisition, frequent secondary loss (regressive evolution)

Range of genome size variation



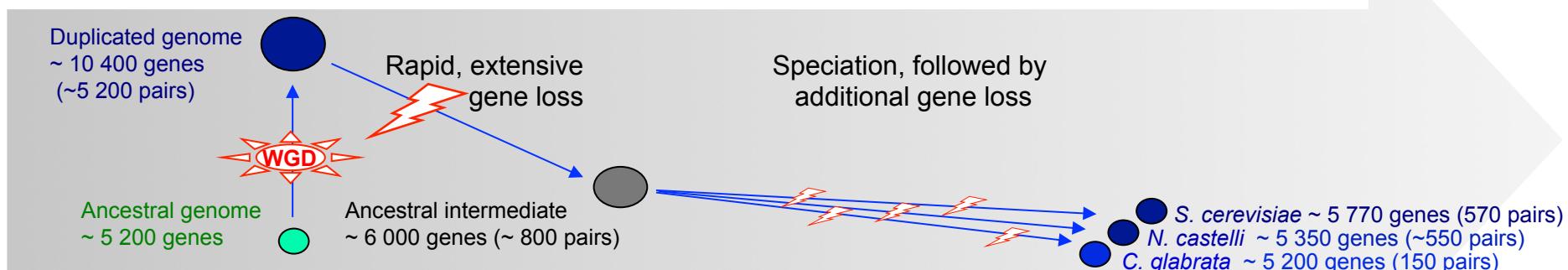
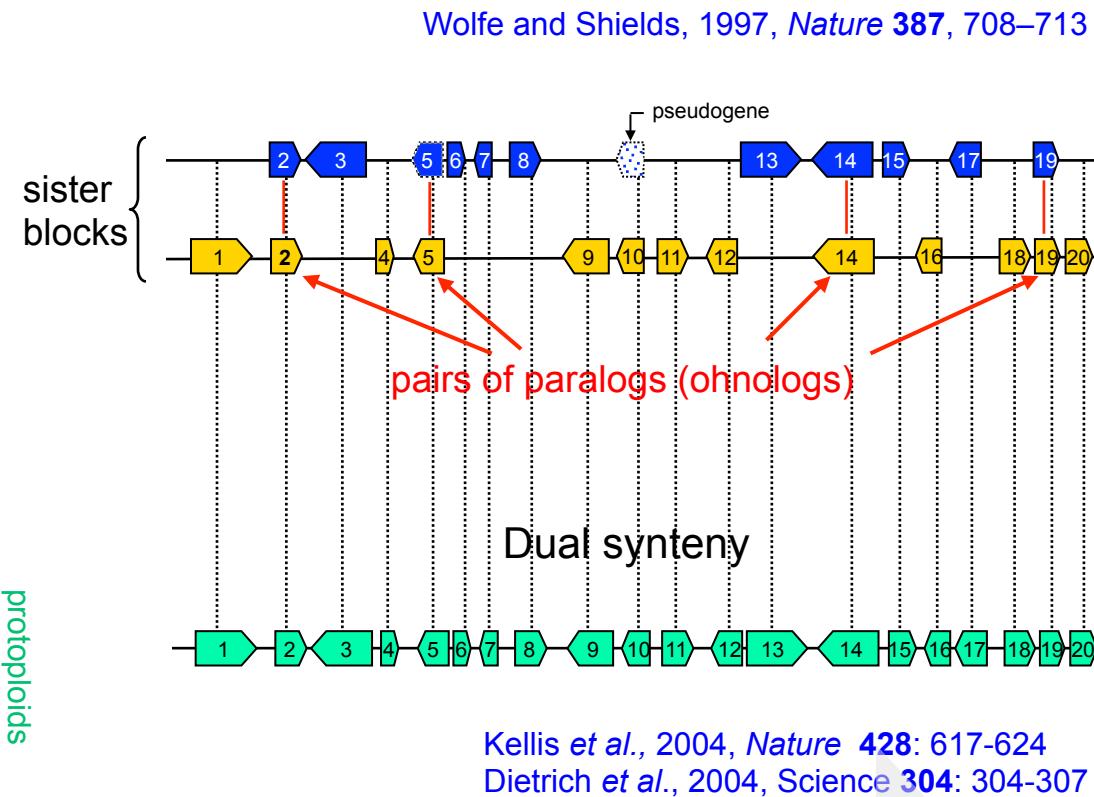
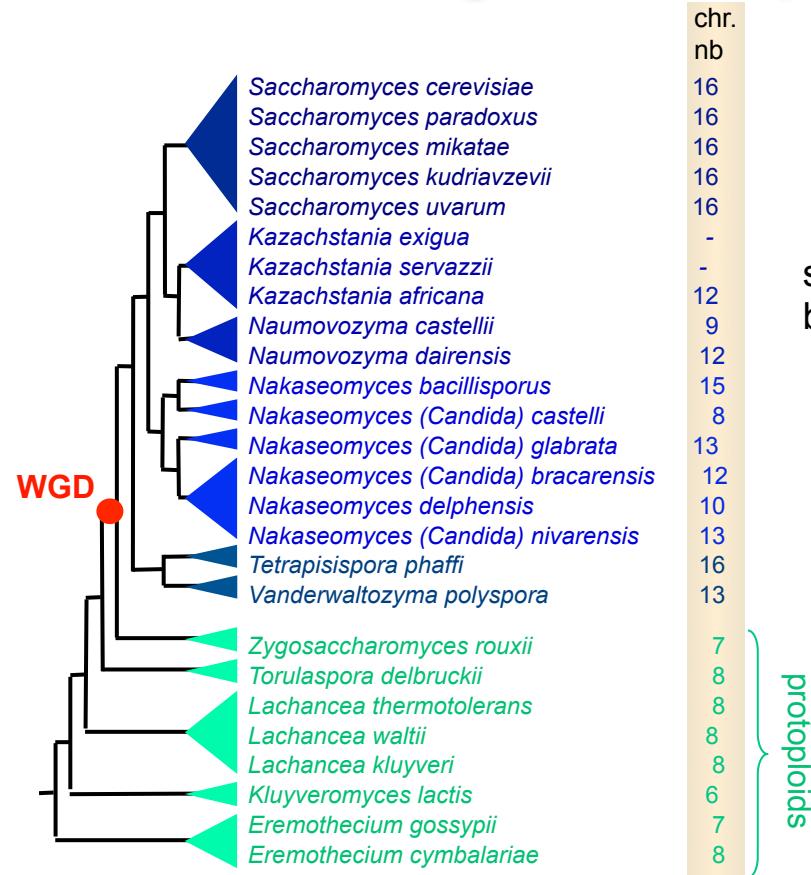


Le monde des eucaryotes, vu par la génomique comparative

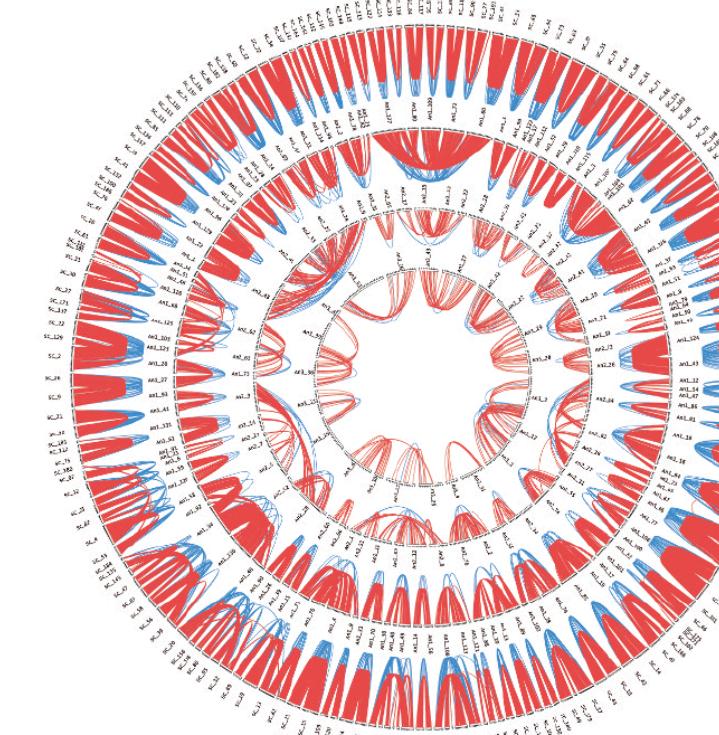
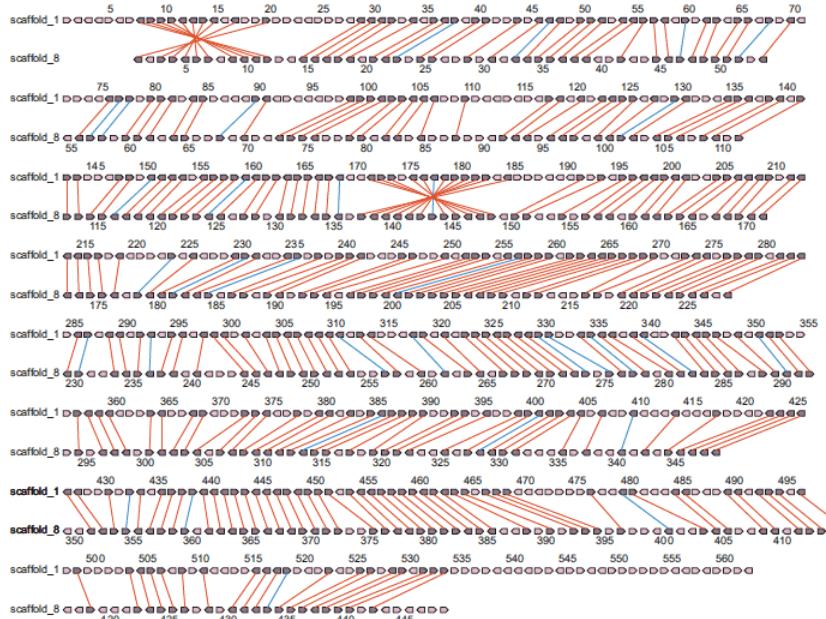
- Les mécanismes moléculaires de l' évolution des génomes eucaryotes

La formation *de novo* de gènes, *frameshifts* et ARNs

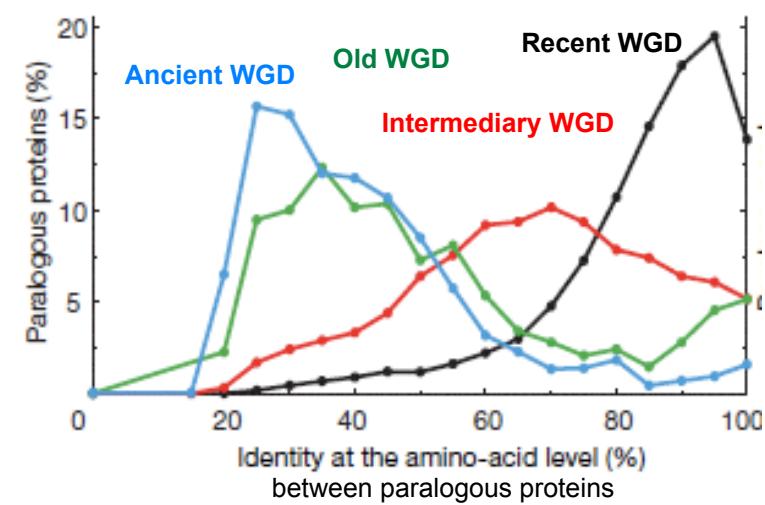
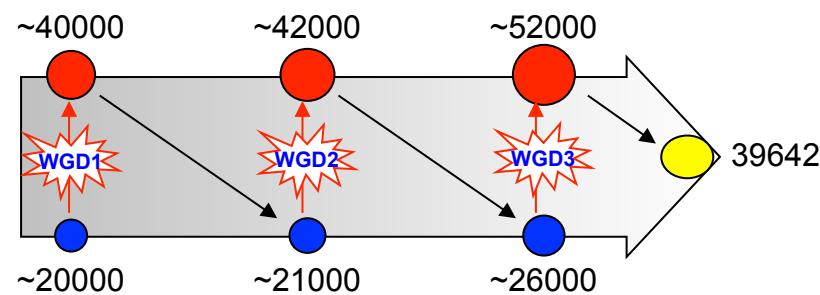
Whole-genome duplication in *Saccharomycetaceae*



Several successive genome duplications in *Paramecium*



Comparison of two scaffolds originating from a common ancestor after a recent WGD in *Paramecium. tetraurelia*

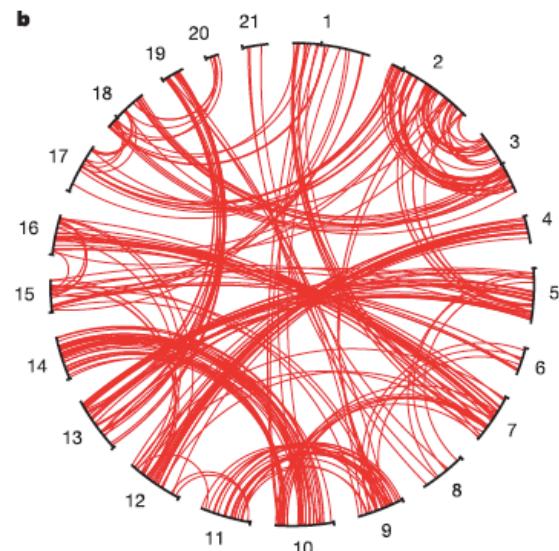


Whole-genome duplication at the origin of *Actinopterygii* fishes

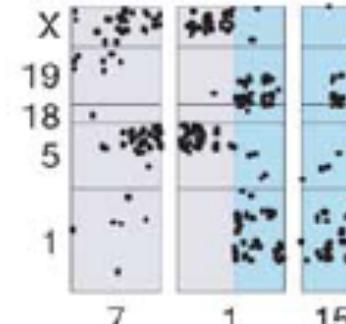
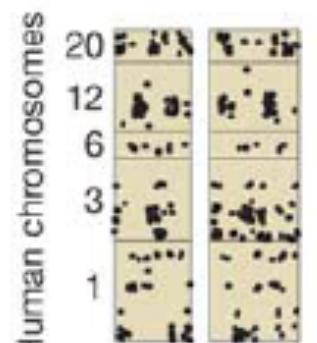
Jaillon et al. 2004 *Nature* 431: 946-957



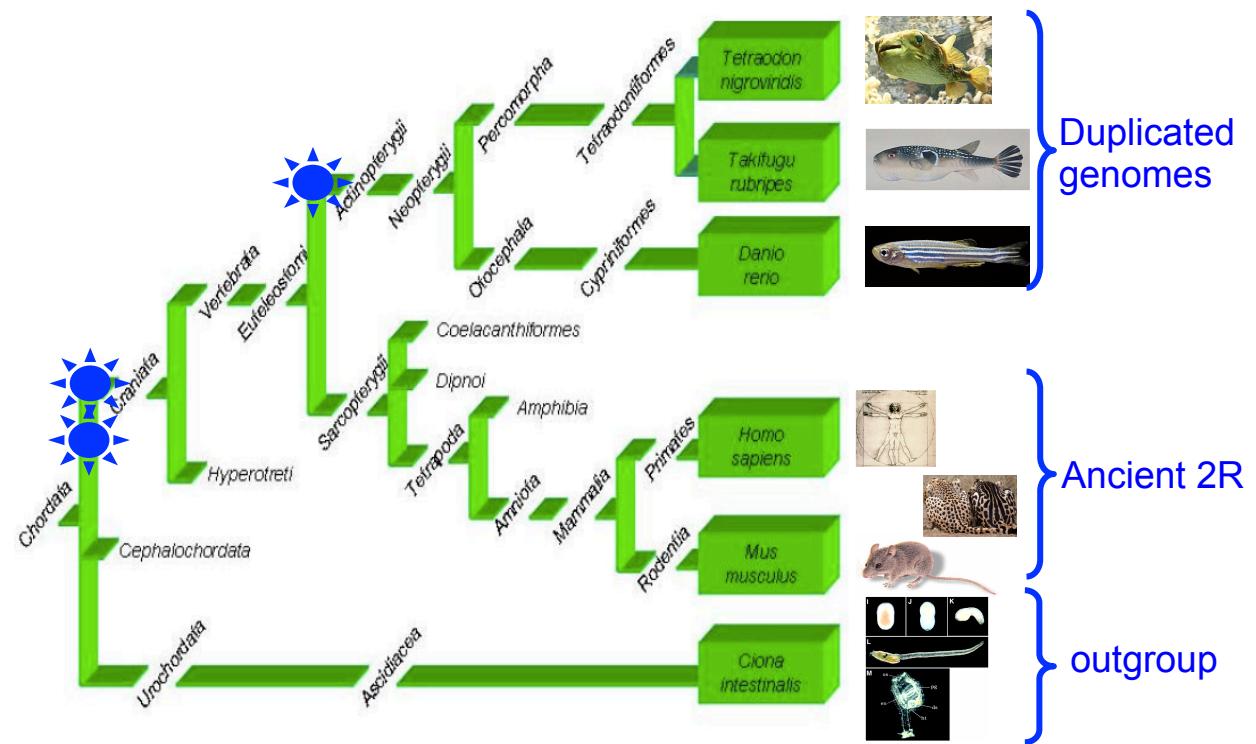
Tetraodon nigroviridis



Sister blocks of synteny



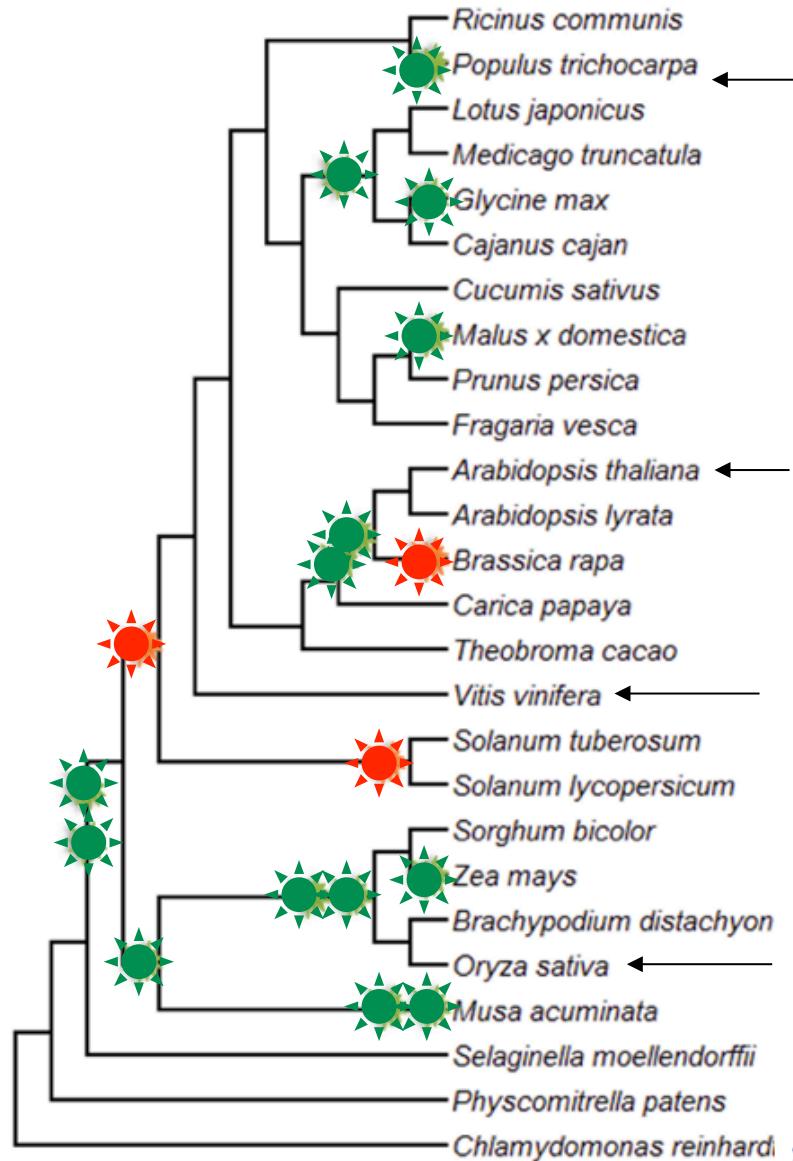
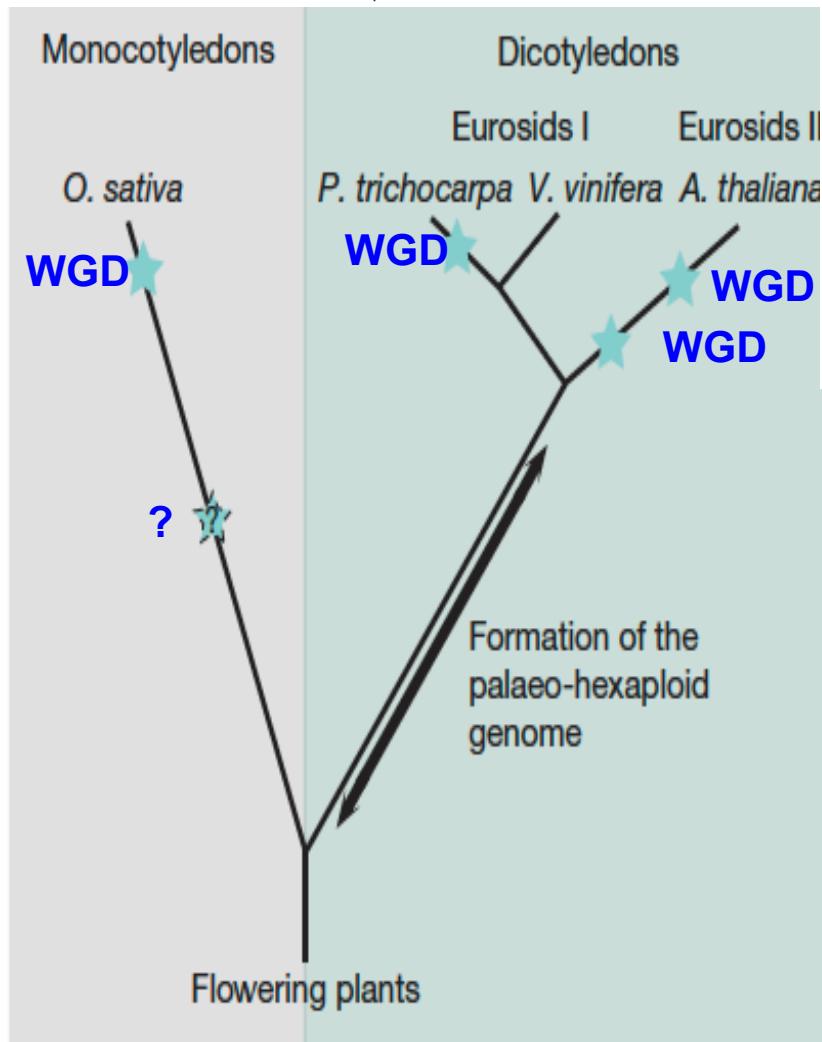
Dual synteny
(Human-*Tetraodon*)



Multiple genome duplications and triplication in Streptophyta

The grapevine genome sequence suggests ancient hexaploidization in major angiosperm phyla

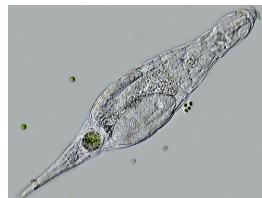
The French-Italian Public Consortium for Grapevine Genome Characterization*



Atypical genomes in metazoa

Genomic evidence for ameiotic evolution in the bdelloid rotifer *Adineta vaga* *Nature* (2013) 503: 453-457

Jean-François Flot^{1,2,3,4,5,6}, Boris Hespeels^{1,2}, Xiang Li^{1,2}, Benjamin Noel³, Irina Arkhipova⁷, Etienne G. J. Danchin^{8,9,10}, Andreas Hejnol¹¹, Bernard Henrissat¹², Romain Koszul¹³, Jean-Marc Aury³, Valérie Barbe³, Roxane-Marie Barthélémy¹⁴, Jens Bast¹⁵, Georgii A. Bazykin^{16,17}, Olivier Chabrol¹⁴, Arnaud Couloux³, Martine Da Rocha^{8,9,10}, Corinne Da Silva³, Eugene Gladyshev⁷, Philippe Gouret¹⁴, Oskar Hallatschek^{6,18}, Bette Hecox-Lea^{7,19}, Karine Labadie³, Benjamin Lejeune^{1,2}, Oliver Piskurek²⁰, Julie Poulaïn³, Fernando Rodriguez⁷, Joseph F. Ryan¹¹, Olga A. Vakhrusheva^{16,17}, Eric Wajnberg^{8,9,10}, Bénédicte Wirth¹⁴, Irina Yushenova⁷, Manolis Kellis²¹, Alexey S. Kondrashov^{16,22}, David B. Mark Welch⁷, Pierre Pontarotti¹⁴, Jean Weissenbach^{3,4,5}, Patrick Wincker^{3,4,5}, Olivier Jaillon^{3,4,5,21*} & Karine Van Doninck^{1,2*}

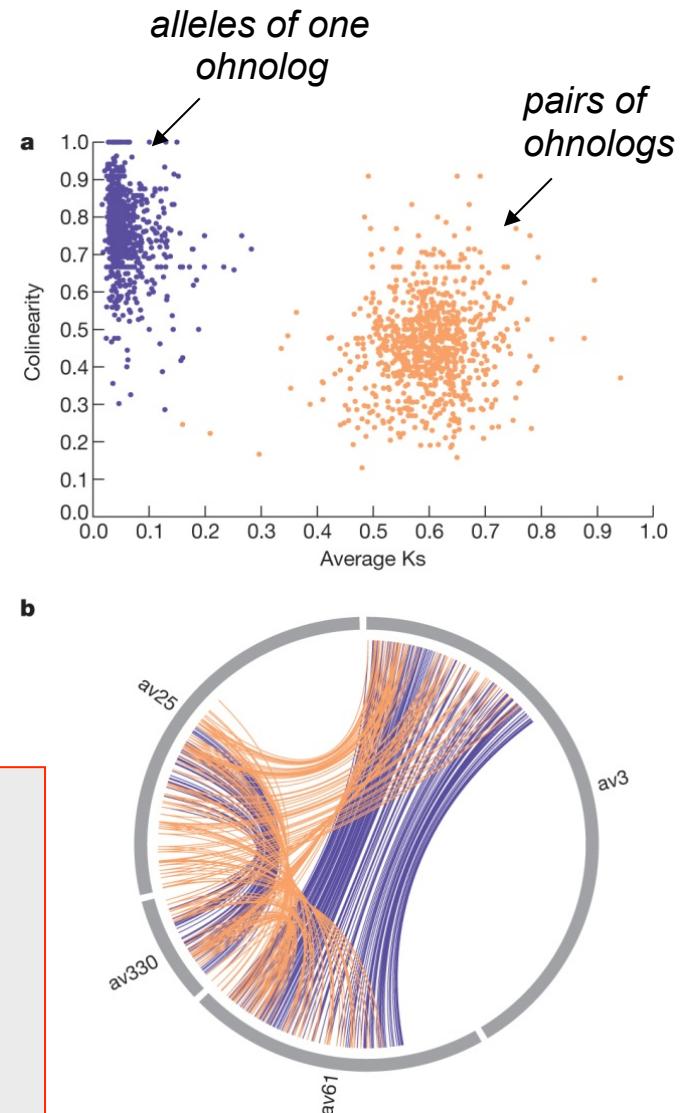


Total genome size 244 Mb
(including 26 Mb homozygous 2x)
49,300 protein-coding genes

Numerous homologous blocks (collinear regions) forming two groups:

- pairs of **ohnologs** (mean 74 % identity) corresponding to an ancient genome duplication
- pairs of **alleles** (mean 96 % identity) for each member of the pair of ohnologs.

Their coexistence in the same genome forms **quartets** that, altogether, cover 40 % of the entire genome (---> a **locally tetraploid** genome)

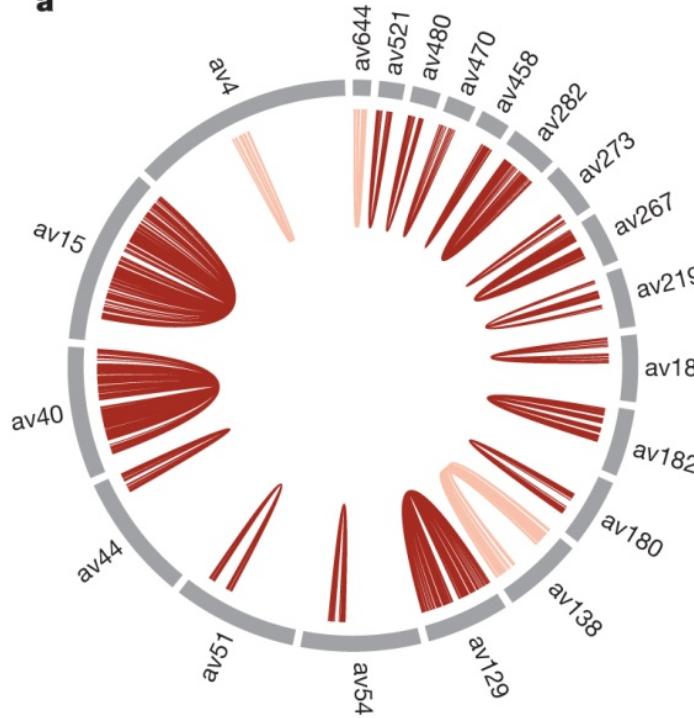


Example of a genomic quartet of 4 scaffolds

Atypical genomes in metazoa

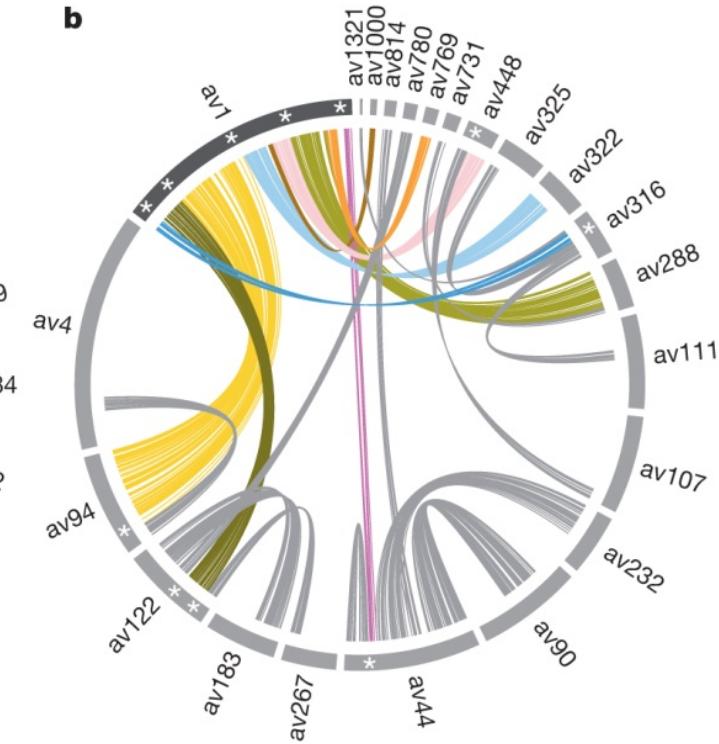
The genomic structure is incompatible with conventional meiosis

a



Allelic pairs are found on
the **same** chromosomes

b

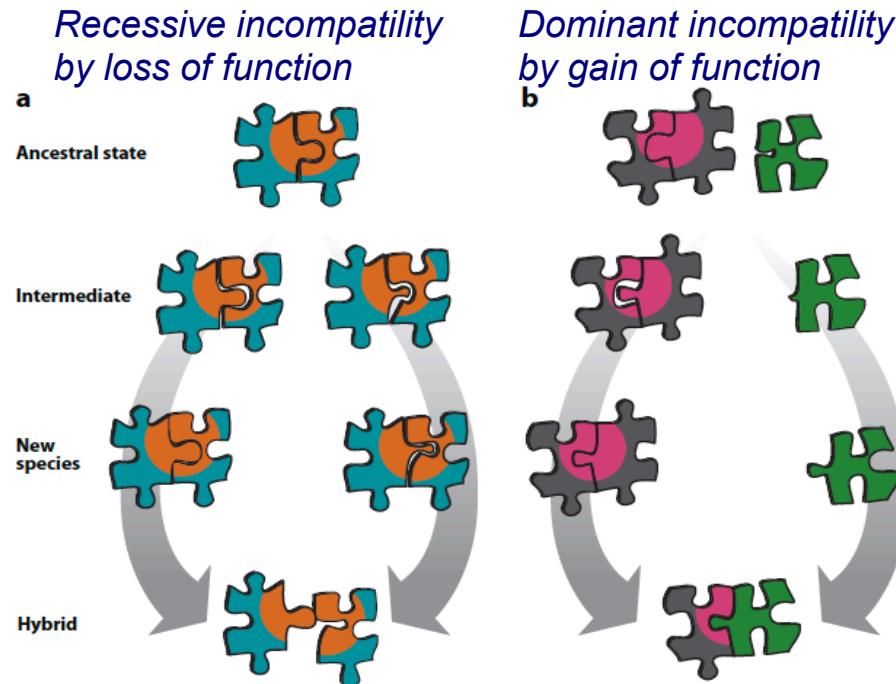


Colinear regions do not extend
to **entire** chromosome scale

Multiple traces of **gene conversion** between gene copies reassort alleles without meiosis

Multiple traces of **horizontal gene acquisition** of non-metazoan origin throughout the genome (8%)

Bateson-Dobzhansky-Müller incompatibility



However: very few examples of DM gene pairs identified

Drosophila: *Lbr - Hmr: heterochromatin turnover*
Zbr: pericentric satellite DNA
OdsH - Ovd: heterochromatin-mediated sex ratio distortion

Xiphophorus: *Xmrk2 deregulated expression*

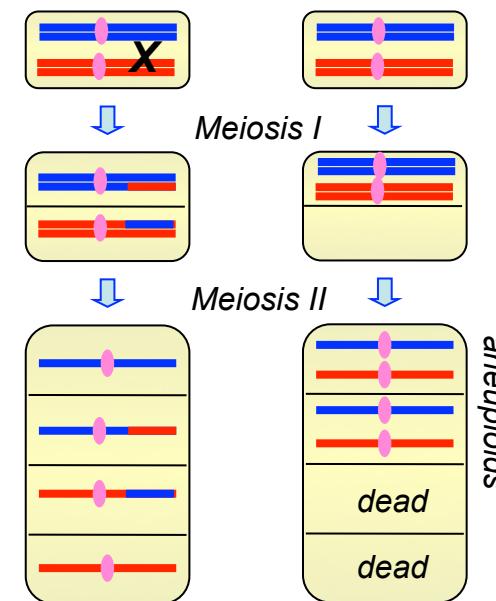
Mus: *Prdm9 meiotic recombination hotspot*
Oryza: *S5: proteolytic enzyme in cell wall*
SaM -SaF: SUMO E3 ligase

Saccharomyces: *nucleo-mitochondrial interactions*
AEP2 (nuc) - Oli1 (mt): lack of Oli1 translation
MRS1 (nuc) - Cox1 (mt): lack of Cox1 intron splicing

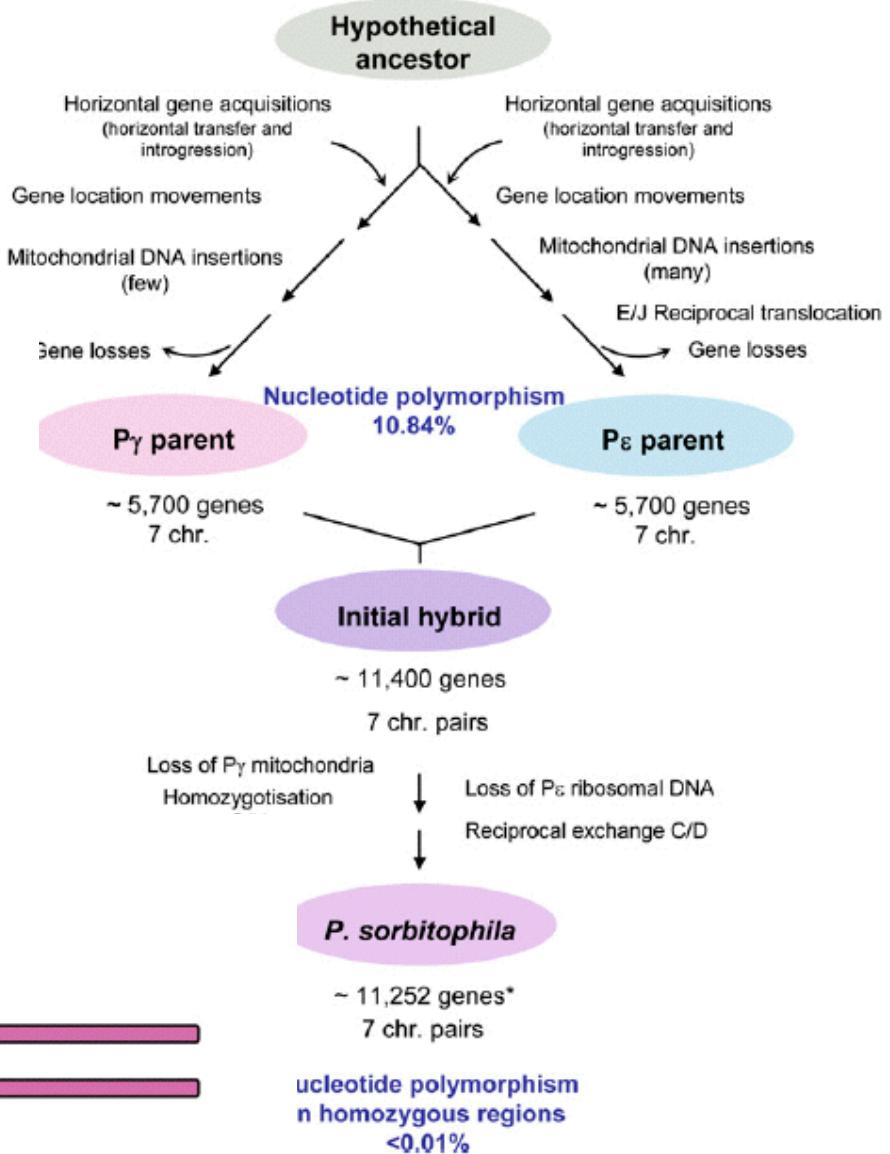
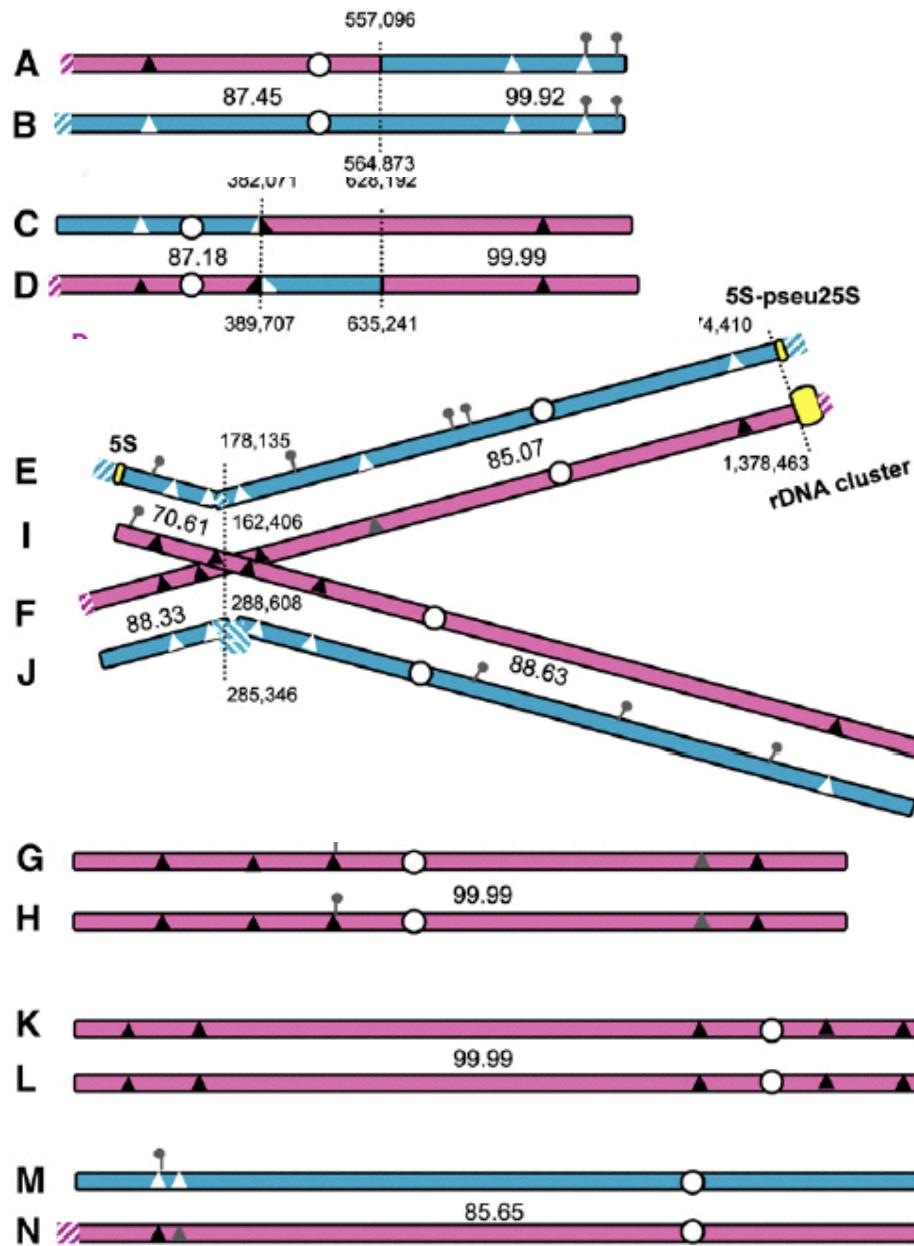
Hybrid meiotic sterility



Sequence divergence --> mismatch repair system --> no COV --> aneuploidy

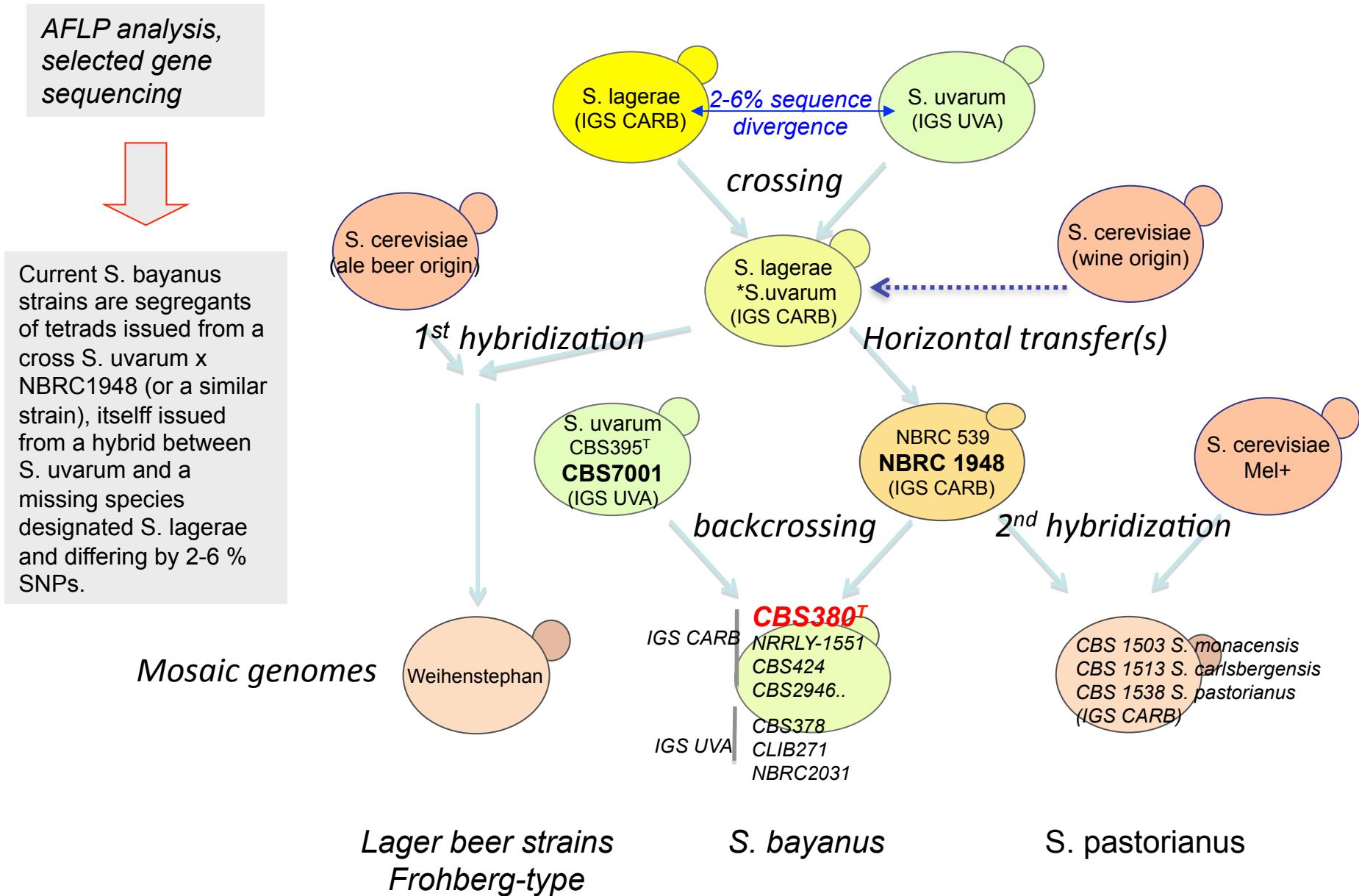


Pichia sorbitophila: a recent hybrid in the process of resolution



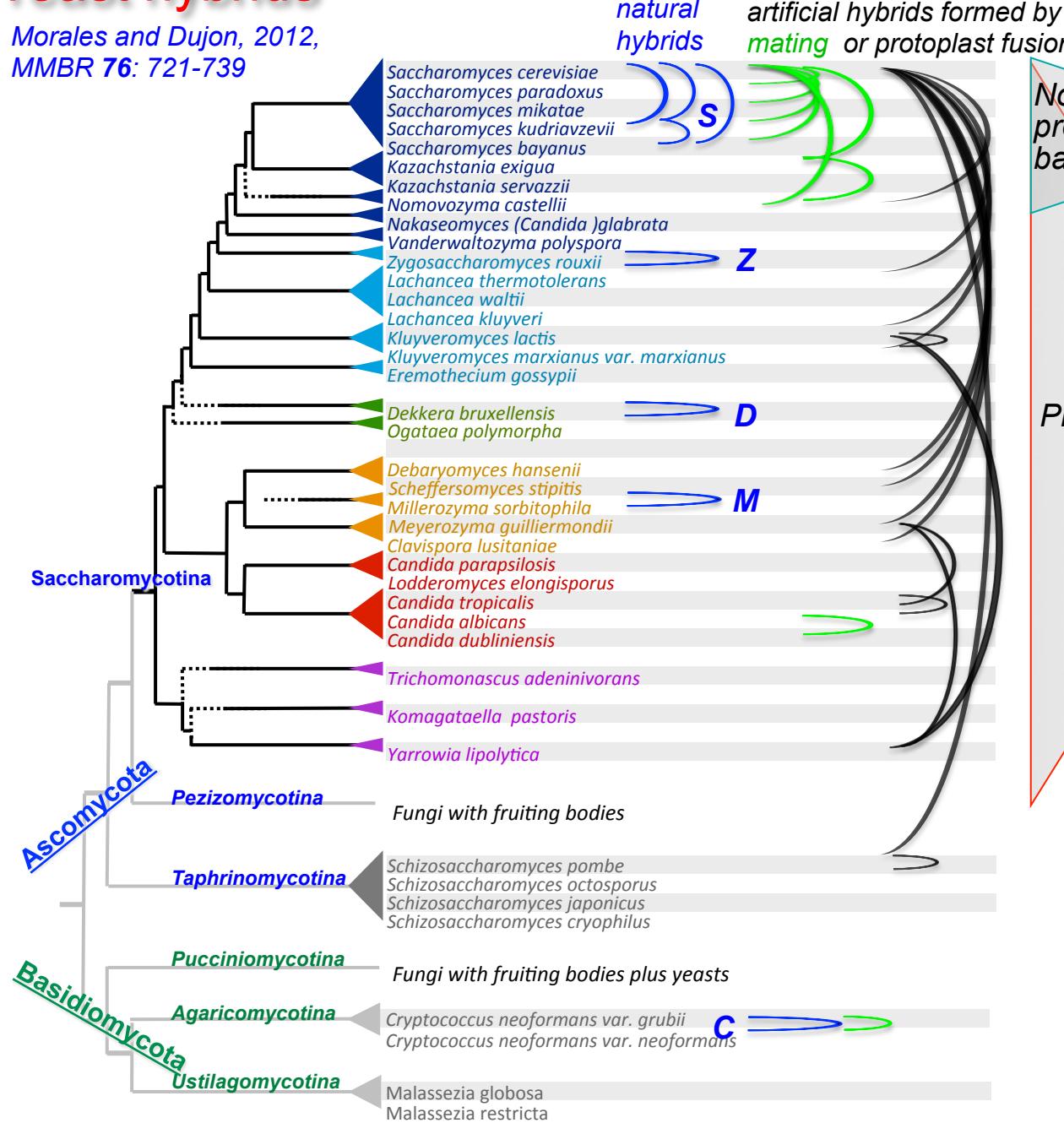
Even type-strains may have a complex hybrid ancestry

Nguyen et al., 2011. PLoS ONE 6: e25821



Yeast hybrids

Morales and Dujon, 2012,
MMBR 76: 721-739



Natural hybrids

S: Saccharomyces
Z: Zygosaccharomyces
D: Dekkera:
M: Millerozyma
C: Cryptococcus

No pre-zygotic barriers

Pre-zygotic barriers

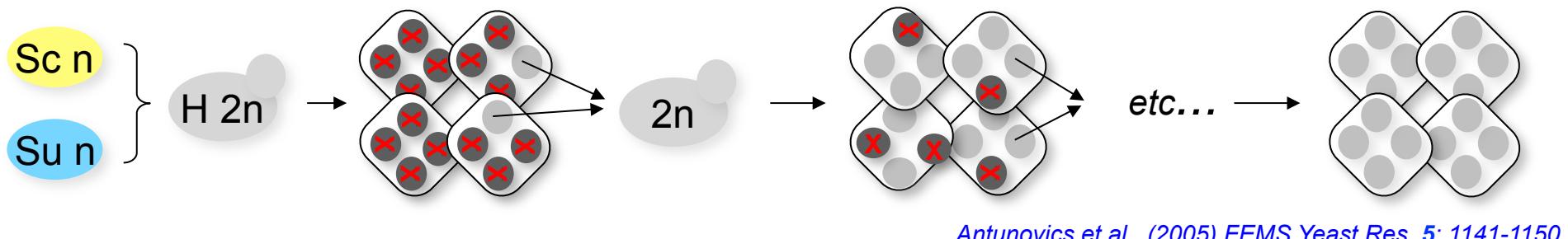
polyethylene, calcium-induced hybridizations between yeasts at considerable evolutionary distances



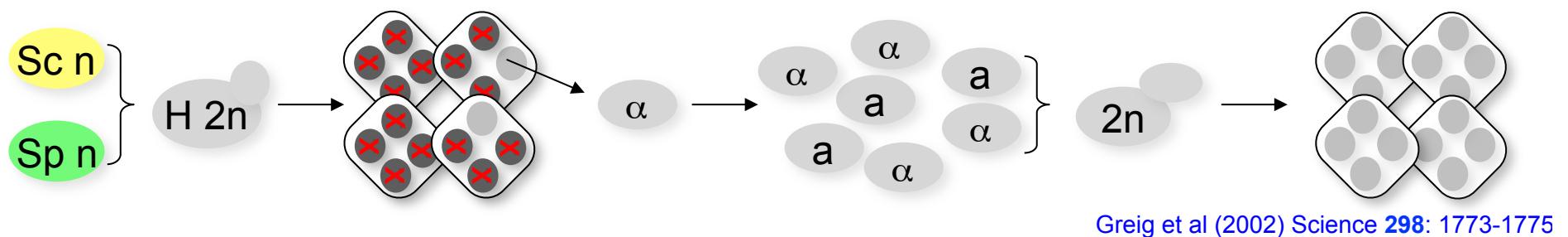
No major dominant Dobzhansky-Müller incompatibilities

Spontaneous processes rapidly restore fertility

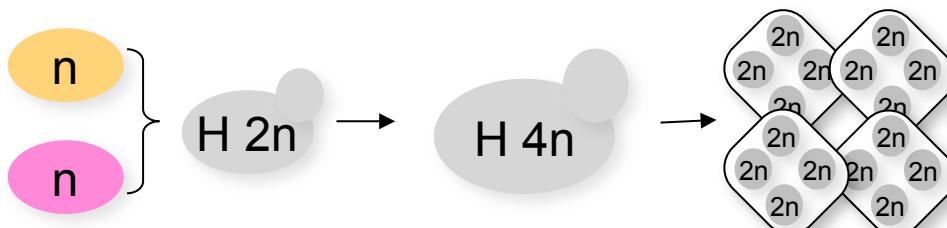
1- Gradual recovery of fertility during successive meiotic cycles



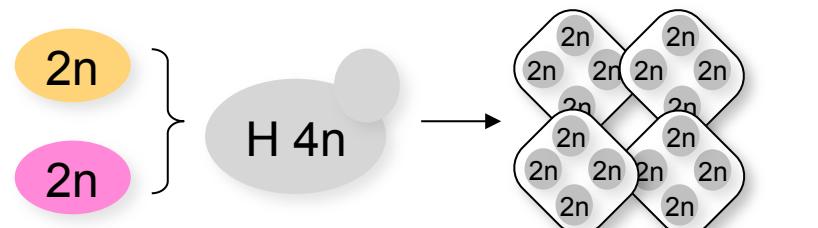
2- Recovery of fertility after mating-type switching

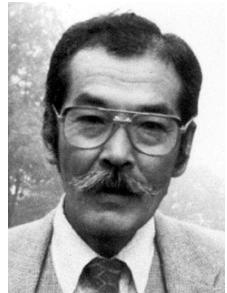


3- Recovery of fertility by endoreplication



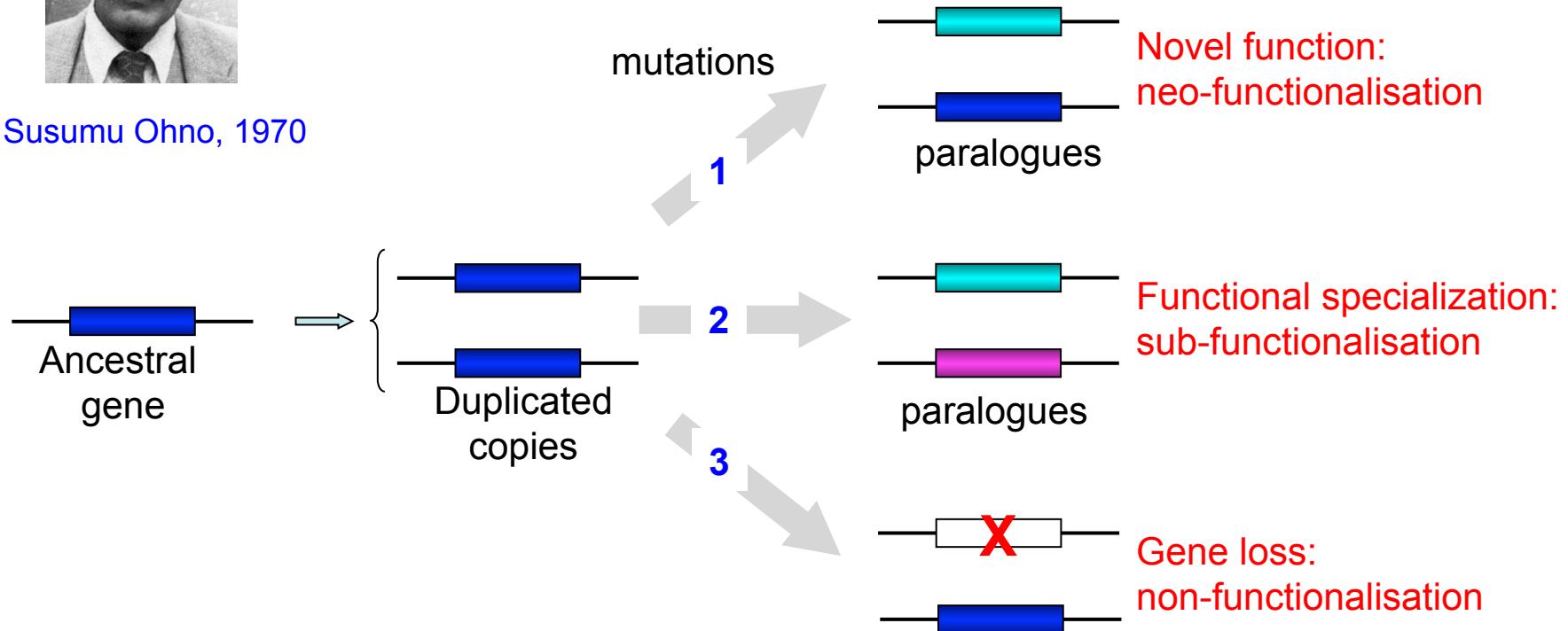
4- Fertility of allotetraploids





Consequences of gene duplications

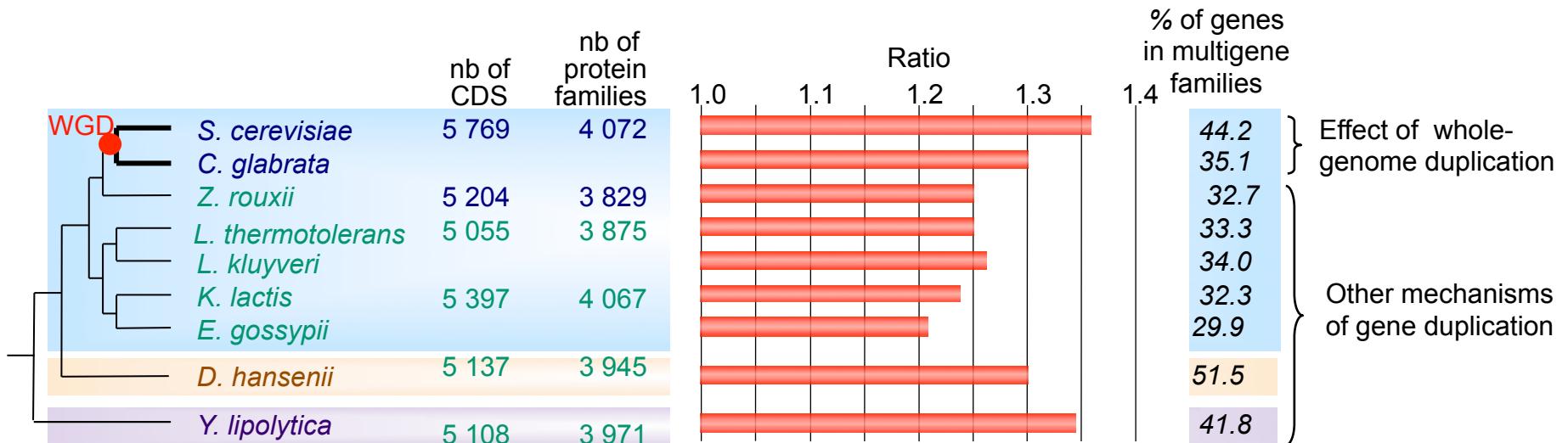
Susumu Ohno, 1970



All genomes show high number of paralogous gene copies resulting from the combination of whole-genome duplications, segmental duplications, tandem gene duplications

A genome is only a snapshot in time of continuous gene duplications and losses.

Duplications and genome redundancy



Single gene duplications:

- in tandem: few arrays in most yeast genomes, more arrays in some species
- ectopic: no direct indication that this mechanism exists
- through RNA intermediate: phenomenon is experimentally demonstrated, but extant retrogenes are very difficult to identify in yeast genomes

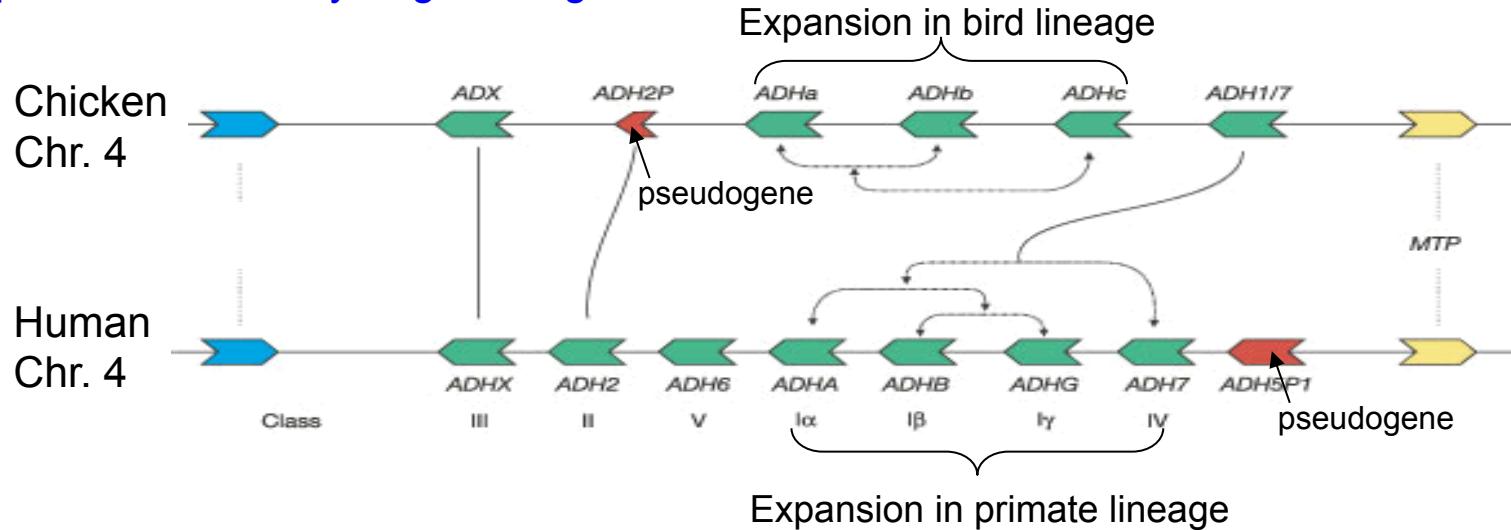
many dispersed paralogs owing to very ancient evolutionary lineages

Segmental duplications:

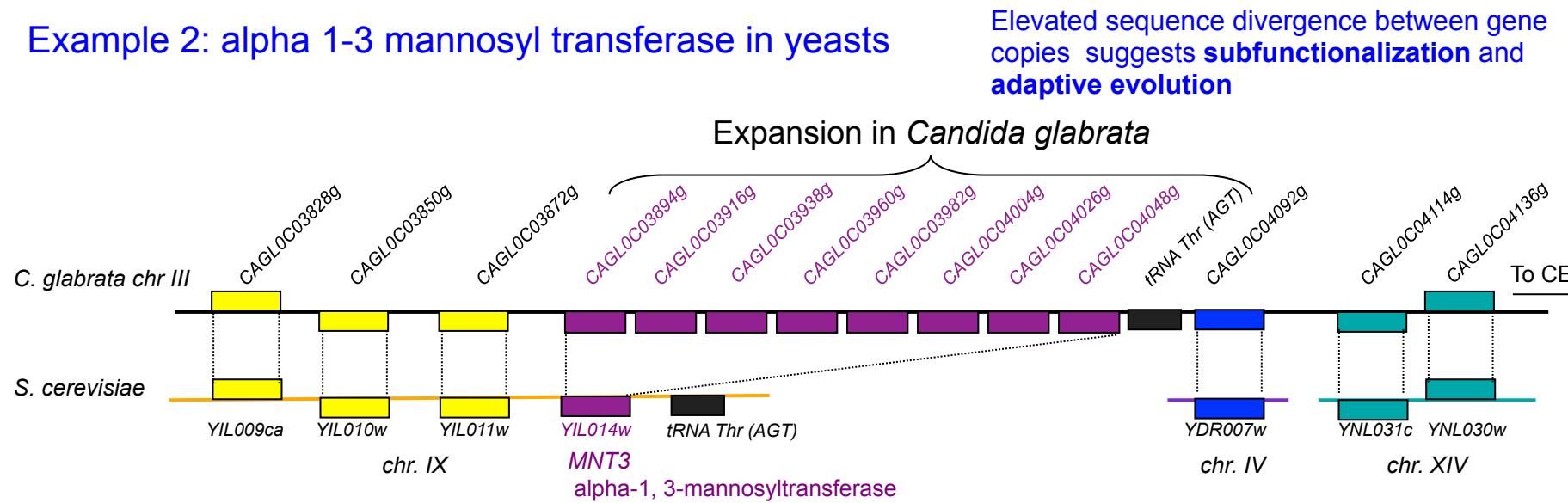
phenomenon is experimentally demonstrated, but examples are very rare in yeast genomes, except in subtelomeric regions

Formation and evolution of tandem gene arrays

Example 1: alcohol dehydrogenase gene cluster in vertebrates

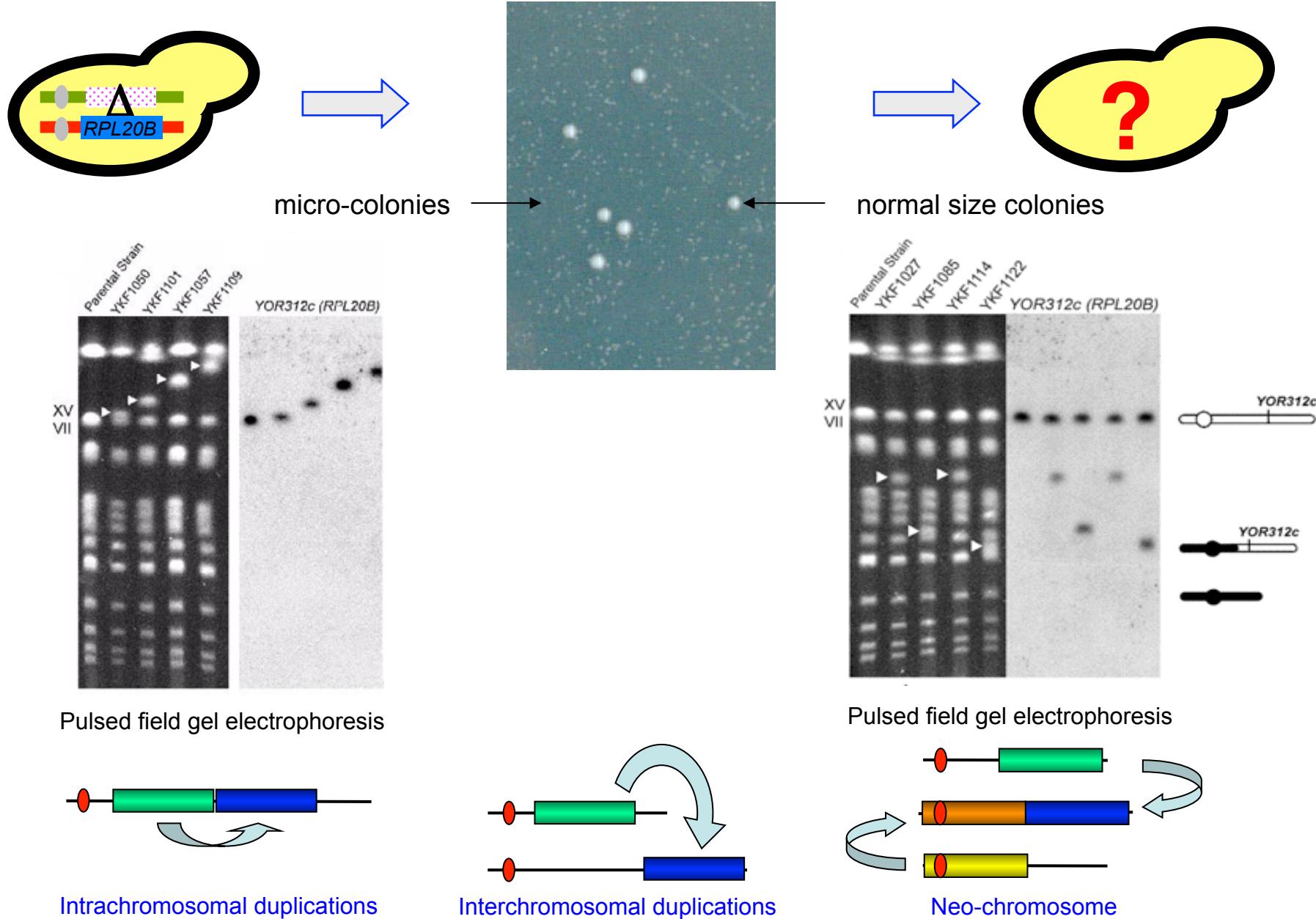


Example 2: alpha 1-3 mannosyl transferase in yeasts



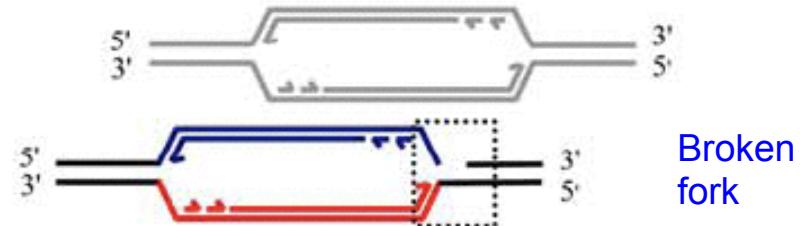
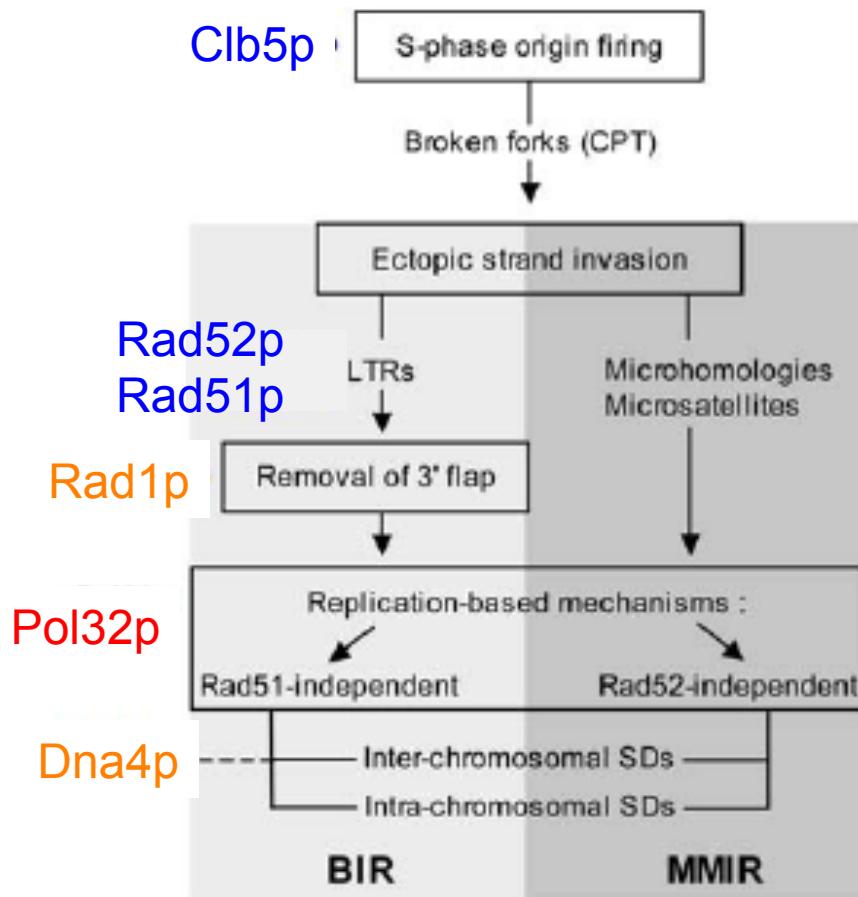
Segmental duplications in *Saccharomyces cerevisiae*

Koszul et al. 2004 *EMBO J.* 23: 234-243

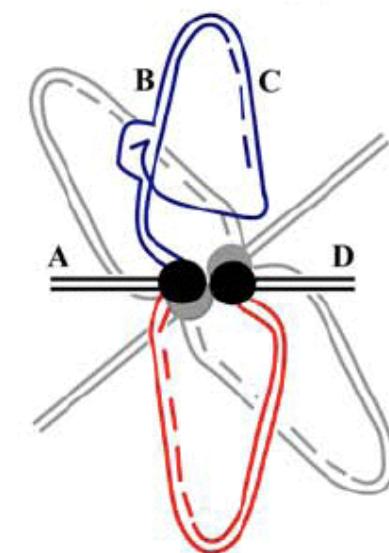


Molecular mechanism of formation of segmental duplications

Payen et al. 2008 *PLoS Genetics* 4: e1000175



Abnormal reinitiation

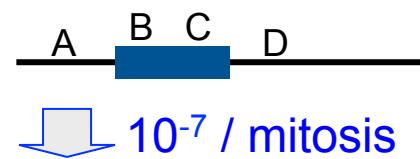


Effect of gene mutation on the formation of segmental duplications.

Red: abolished > essential protein

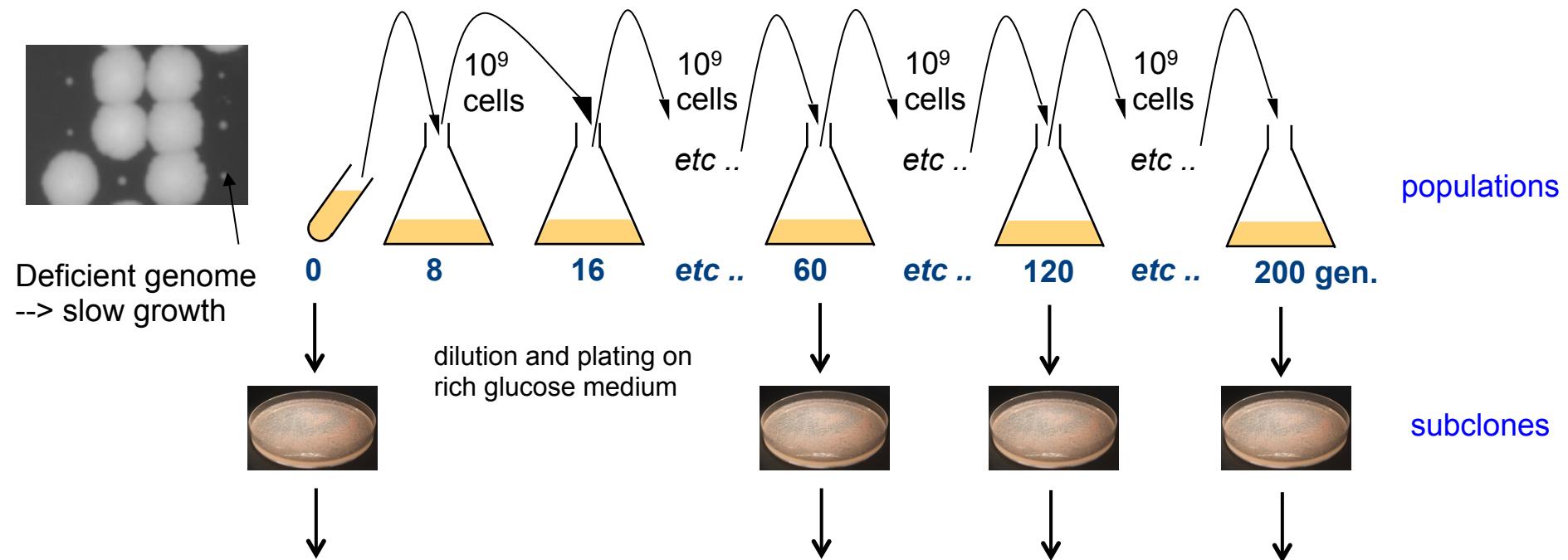
Orange: reduced > protein involved

Bleu: increased > control protein



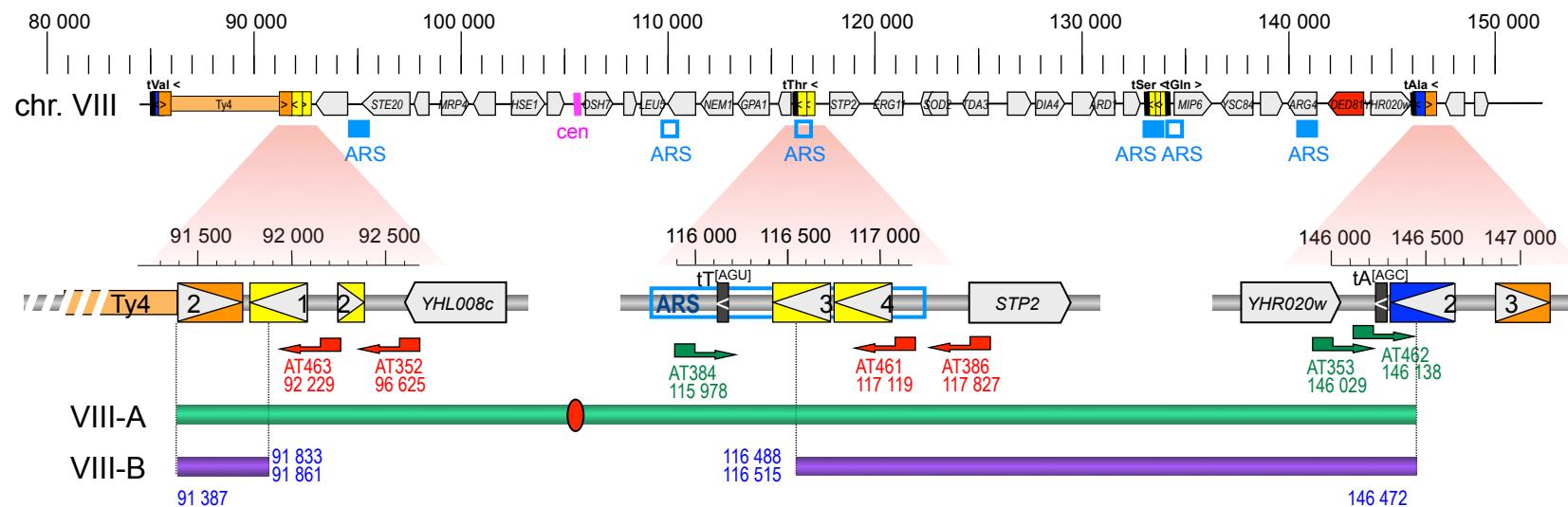
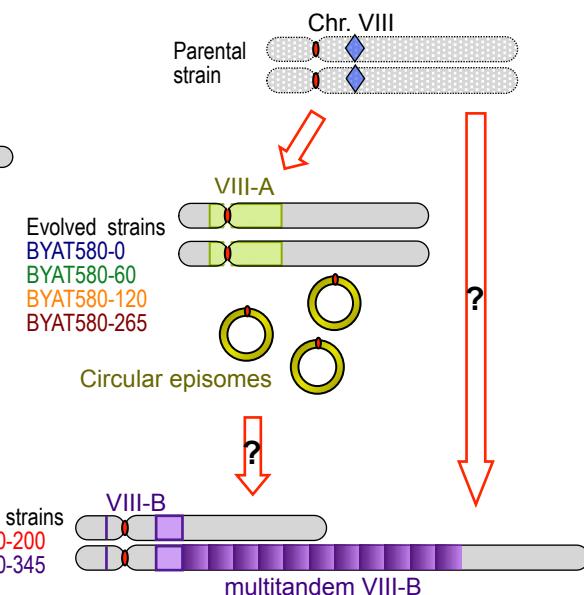
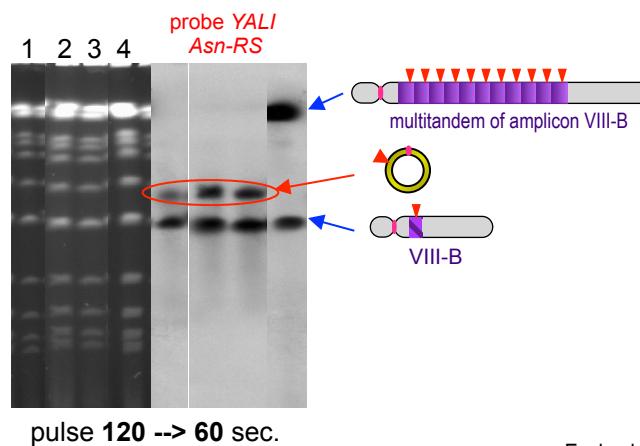
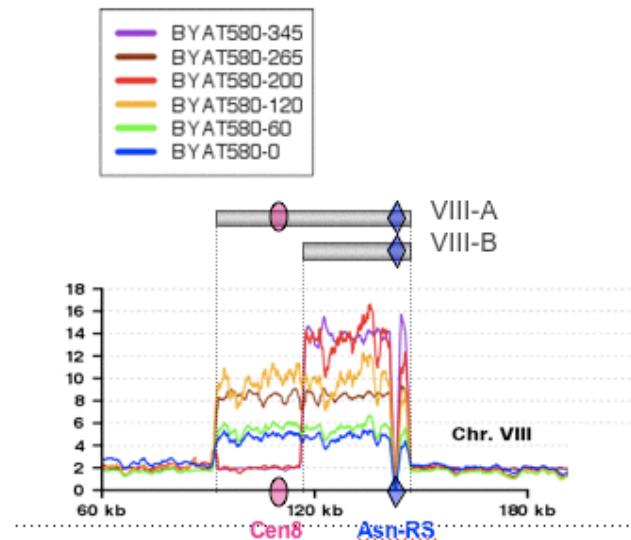
Experimental evolution by serial transfer of populations

- two-liter batch cultures in rich glucose medium at 30°C with aeration
- serial transfers with population bottlenecks of 10^9 cells



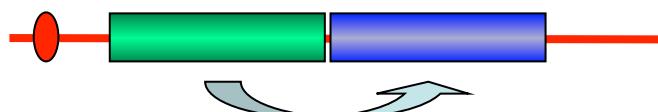
Fitness measurement (growth rate)

Whole genome deep sequencing (Solexa HiSeq2000)

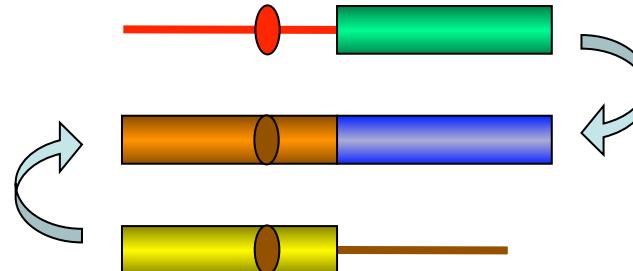


Types of chromosomal dynamics

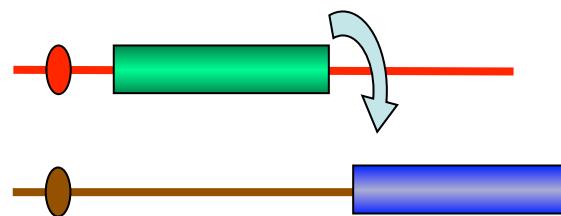
(a) Intrachromosomal



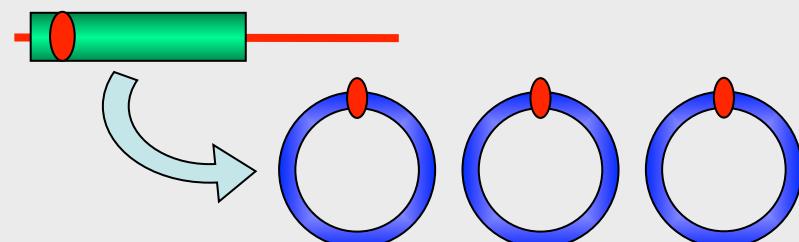
(d) Neo-chromosome



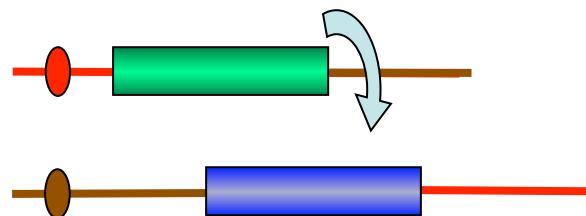
(b) Interchromosomal (subtelomeric addition)



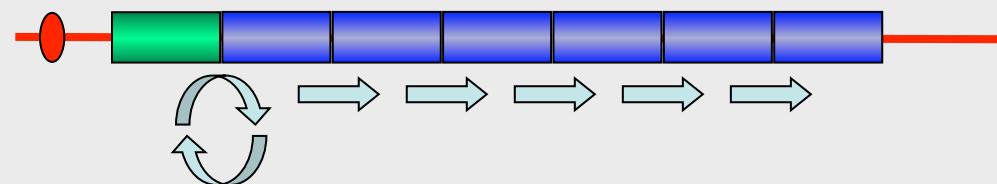
(e) Episome



(c) Interchromosomal (reciprocal translocation)

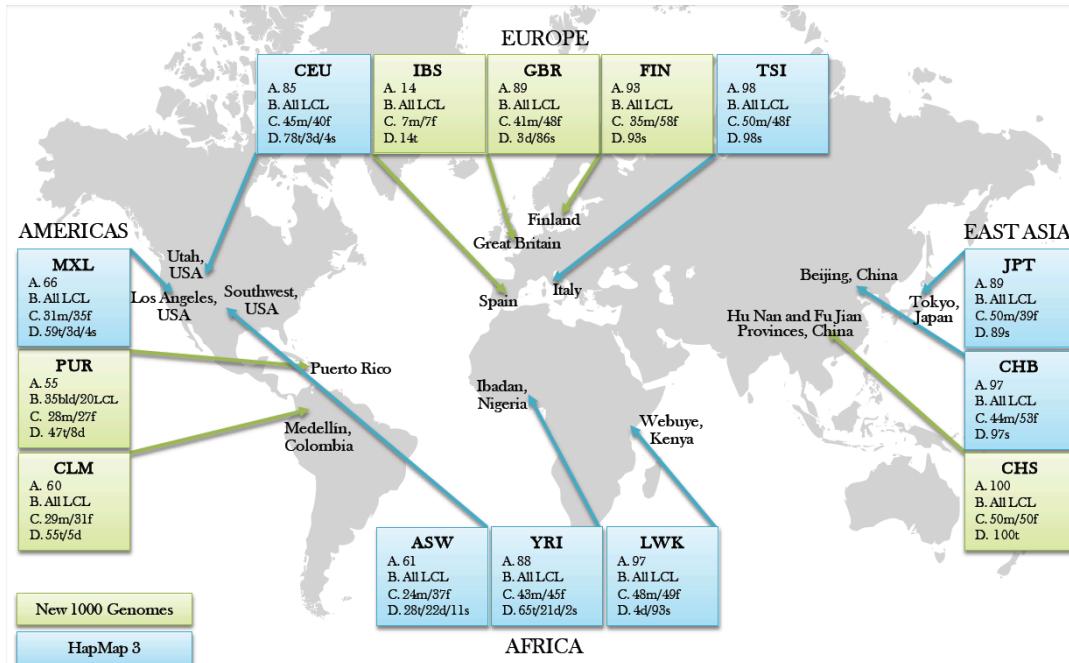


(f) Multitandems



Polymorphic variation of the human genome

The 1000 Genomes Project Consortium (2012) *Nature* **491**: 56-65



1 092 adults from 14 populations

Sequencing method

genome-wide low coverage
high coverage of coding exons
a few parents-child trios

Results

38 millions of SNPs
1.4 millions of indels
14 000 large délétions

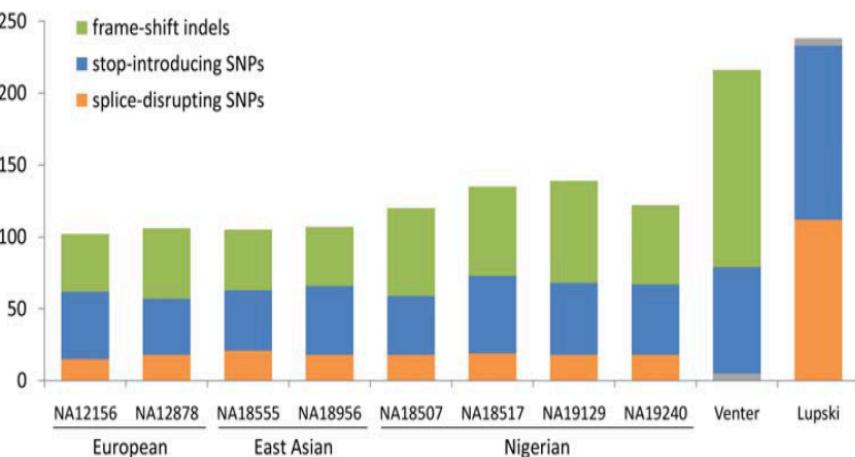
Compared to the reference genome, each individual genome carries:

3 millions SNPs (~ 40 % homozygote)

400 000 indels (~ 40 % homozygote)

Tens of deletions and duplications of average size 8 kb, totalling ~ **400 kb** (0,01 % of genome)

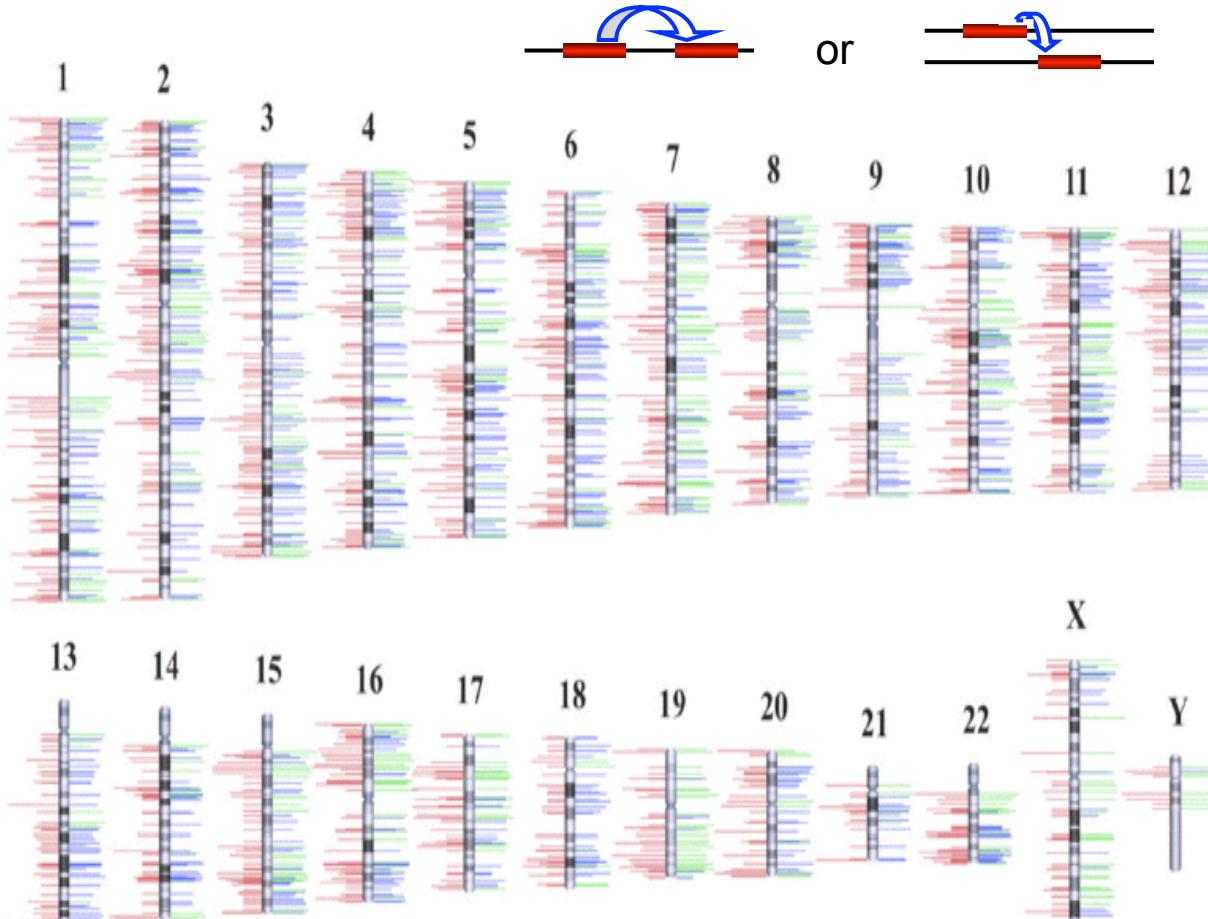
In total > **100 loss of function** per genome (including **30 homozygote**)



MacArthur and Tyler-Smith (2010) *Human Mol. Genetics* **19**: R125-130

Duplications and deletions in normal individuals

Common sites of segmental duplications



Study of 270 normal individuals
(Europa, Africa, Asia)

Number of cases	Size (kb)
- 1	<10
- 10	100
- 100	1000

- 1 447 sites identified (~ 12 % of genome)
- each individual carry duplications covering ~ 5% of genome



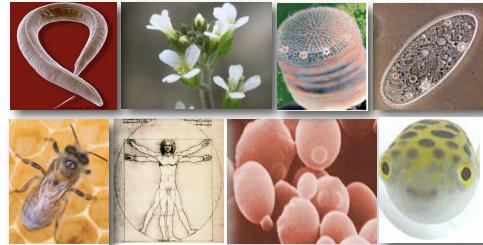
From parent-child trios: ca. 550-750 kb missing at each generation

(Conrad et al., 2006 *Nature Genetics* 38, 75-81;
McCarroll et al., 2006 *Nature Genetics* 38, 86- 92; Hinds et al., 2006 *Nature Genetics* 38, 82- 85)

One new gene created on average for every 21 newborn babies.

60 new genes appeared in the human genome since its separation from chimpanzee. Most new genes are highly expressed in cerebral cortex and testicles.

(Wu et al., 2011, *PLoS Genetics* 7: e1002379)



Le monde des eucaryotes, vu par la génomique comparative

Les mécanismes moléculaires de l' évolution des génomes eucaryotes

- La formation *de novo* de gènes, *frameshifts* et ARNs

Comment naissent les gènes ?

- Réassemblage d' éléments préexistants (au niveau ADN ou au niveau ARN)
- Création *de novo* de séquences codantes ou d' ARN non codants

Mécanismes disponibles au niveau de l' ADN

duplications - divergence de séquence (tandems, segments, génome entier)

fusions / fissions (au niveau des jonctions de segments chromosomiques)

frameshifts

création de nouveaux promoteurs et de nouvelles jonctions intron-exons (divergence de séquence)

transcription-traduction antisens

existence de protogènes

Mécanismes disponibles au niveau de l' ARN

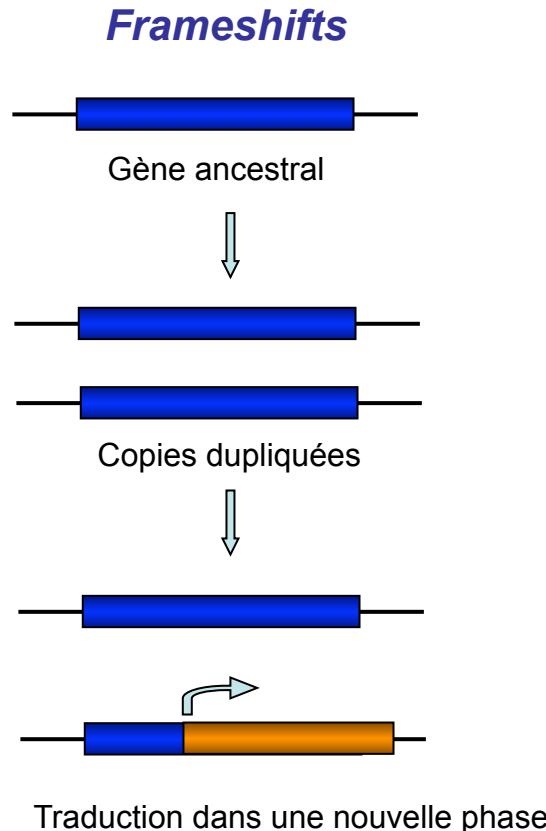
exon shuffling

rétrogènes

exonisation de séquences non-codantes

capture de promoteurs existants et création d' exon 5' non codants

Formation de nouveaux gènes par *frameshifts*



Chez la bactérie *Flavobacterium sp.* (portant un plasmide K172), une **nouvelle fonction enzymatique capable de dégrader le nylon** est apparue spontanément au début des années 1980.

Il s'agit d'une **protéine de 392 acides-aminés** montrant une **activité hydrolase** sur les oligomères linéaires d'acide 6-amino hexanoïque (nylon).

Après examen de la séquence, on trouve que cet enzyme est produit par la **traduction hors de phase** d'une séquence préexistante codant une protéine de 427 acides-aminés **riche en arginine**.

(Ohno 1984 *PNAS* 81: 2421-2425).

Formation de nouveaux gènes par *frameshifts* (génome humain)

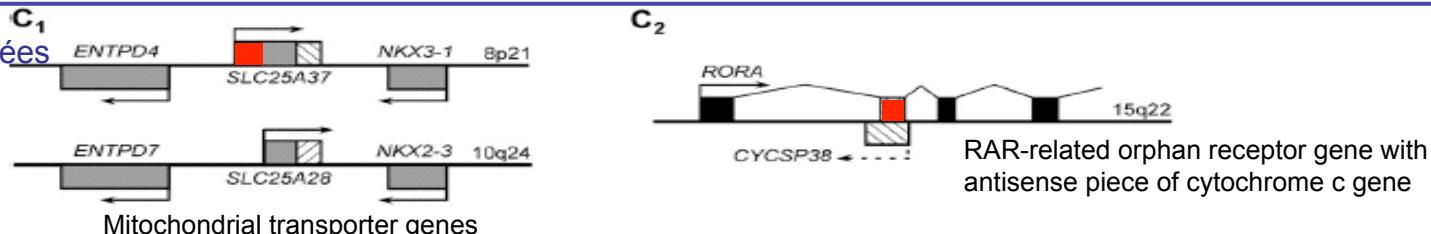
Okamura et al., (2006), *Genomics* 88: 690-697

Examen systématique de nouvelles séquences capables de coder des protéines dans le génome humain par *frameshifts* de gènes ancestraux. Examen de 23 052 séquences de mRNA (RefSeq après curation).

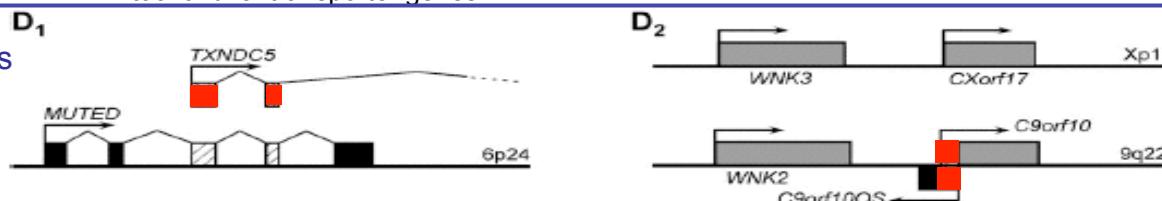
Tandems



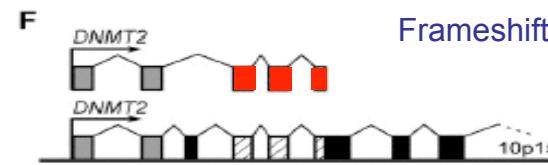
Duplications dispersées



Gènes chevauchants



Frameshifts internes

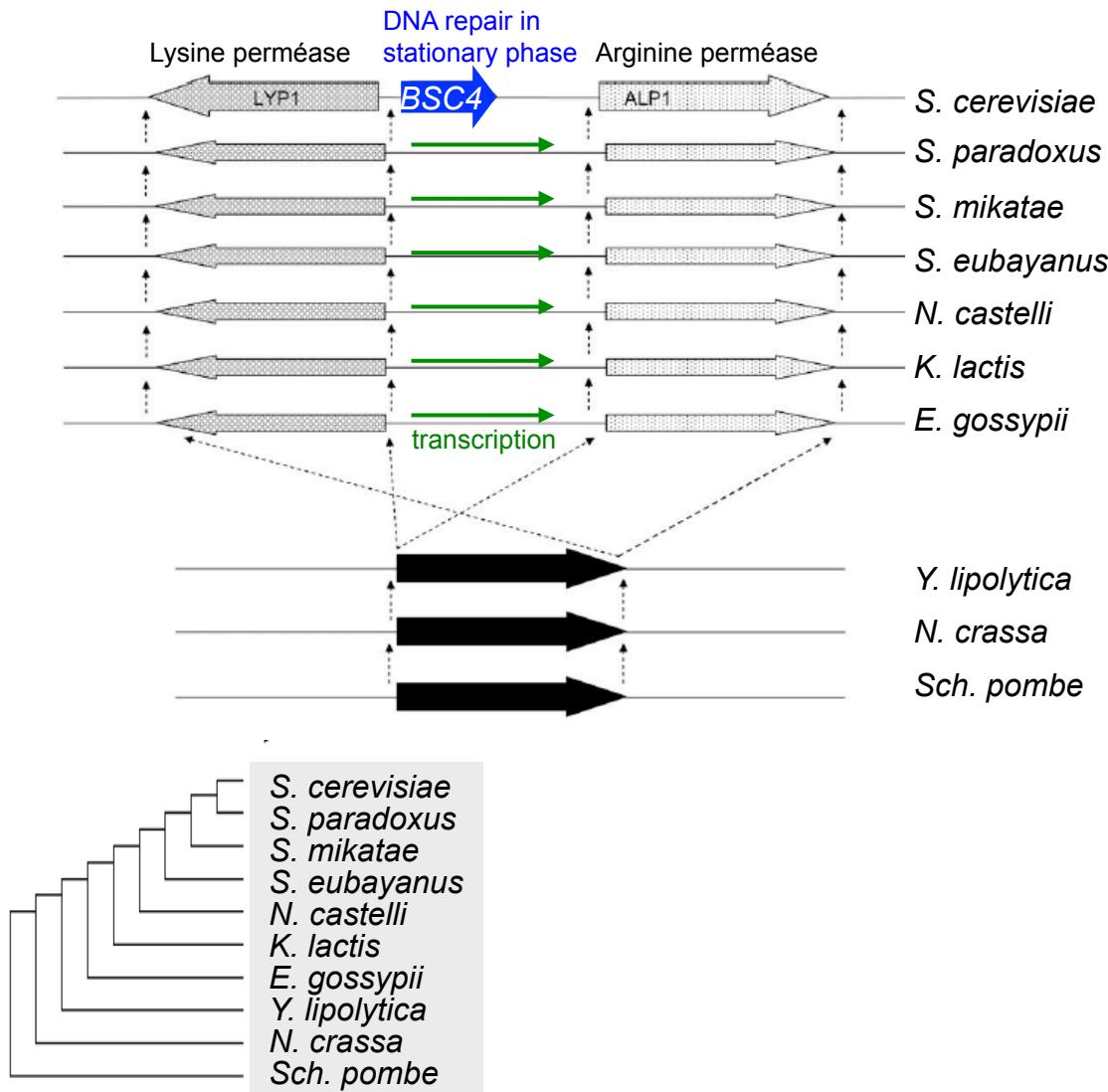


➤ 470 *frameshifts* (17 - 397 codons, moyenne 56 codons) parmi paires de paralogues examinées

Autres données, revue de Raes and Van de Peer 2005 *Trends in Genetics* 21: 428-431

Création de novo de gène dans une zone transcrise

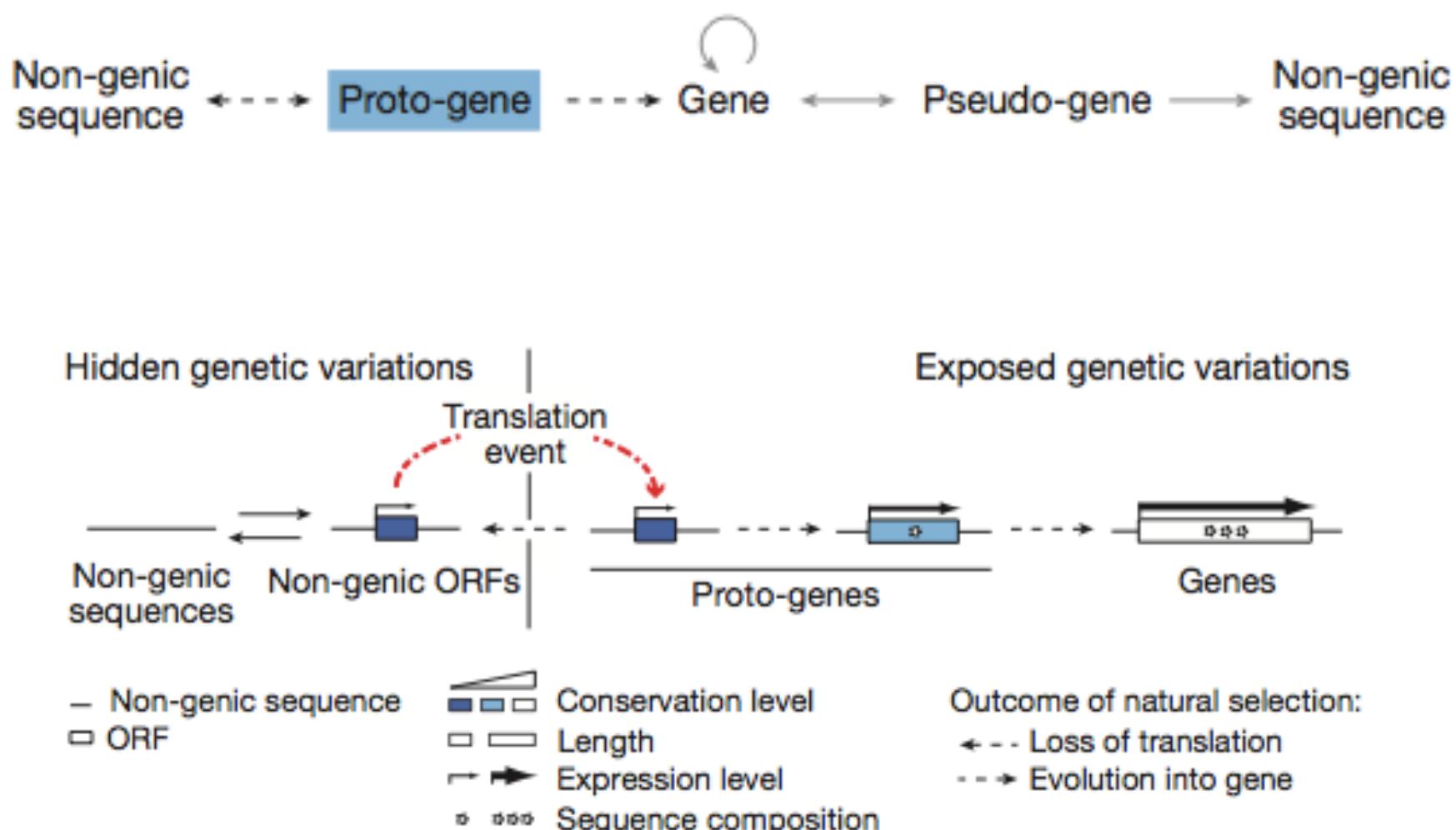
Cai et al., 2008 Genetics 179: 487-496



Scer	AATGTCIAATTGTGCCTACGGCAAGAGTAACAAAAAAACAAAAAAACTGCATAAC[50]
Spar	GTTGTCITGTAATTCTACGGCAAAAGTAAACAAAAAAACTGTATTGCATAAC[50]
Smik	GTGCCCTTAATTCTCGCCCAAGA-CACACAGAAAACCAAAATAATAAT[50]
Sbay	GIACTTGAGCT-TACGCCAGAAAAGAAAAGTTAGAAAAGTGCATGCA[50]
Scer	AAGCAAGT---TTTATACAATACACATTA-----T[100]
Spar	GAGGAATT---TATATACAATACACATAG-----T[100]
Smik	AAACAAAAA---CATATCCGATAACATCGGAG-----T[100]
Sbay	CAGCAAGAGGAGTATATACAGTGGAGACTGGAGAGCCTAGACACACTGT[100]
Scer	AAAA-----ATTCTACTCGGGTTCGGAGCTCCATTGCCATT[150]
Spar	AAAG-----ACTTCGCTCTGATGTC-CGAACCTGCCATTGTCATT[150]
Smik	AAAGTTATA-TCTACTTCACTECACTGCCGCTGACTGT-GTTCGCCATC[150]
Sbay	ATAGGCCACAAAGTCGTTCCGGCTECGTTGCCCGGAATTACATTGCAATT[150]
Scer	GGAGAAAGCCCTTA-TCTGGAGTGGACTGCCCTACAGGTGTTTCAGG[200]
Spar	GGAGAAAATCCCTTA-TCTGGAGTGGACTGCCCTGCAAGTTATTCAGA[200]
Smik	GAAGRAAAATCTTIA-TCTGGATGGAGTACCTGCAAGG-----[200]
Sbay	GGAGAGCCCTGCTAGTGGAGTACAGTATATGICAGC-----[200]
Scer	AAAGACATGGTTACAAAAAAAG-ACGACATTCGCC---AACTTATCACT[250]
Spar	GAAAGCTGGTTACAAAAAGGGACCAGATTCGCC---TAGCTTACAACTC[250]
Smik	ATGCTGTGGTTACAAAMAGGACCAATTACACAGTAACGTATGTTTC[250]
Sbay	GAAGACATGGTTACAAAAAAAGAACCATGCTCGTTACTCACTGTATGCTT[250]
Scer	GCTTGAACCACTTTTATGCCAGGCCCTAACGCCGGAC--TC-AAAA[300]
Spar	GCTTGAATCA-TCTTATGCCAGACCTTCAACGCCGGAC--CCCAAAA[300]
Smik	GCTATAAGCA-TTTTATGTCAGATCCTAACGCCGGAC--CCCAAAA[300]
Sbay	GACTGCTTT-TTTTATGCCAGACCTTAAAGGCCGGCAACCCCCAAA[300]
Scer	ACATACATAC---TGTGGCGCACGGCAGTTTTGG---CGCTATGACACCC[350]
Spar	ACATAATGCTGAGTCACCATGGTCTGGGGCTG-TCGCTGTCGCGCTG[350]
Smik	ACGT---TGC-AGATATTTGA-GTCGCACATCGATAGCTACTACATTA[350]
Sbay	AAACAATATAACGCACACTCGCTTGTGCCGAAGCAAGGCCGTACGCTG[350]
Scer	TTTCCCCAAGAATAT-CGCATTTAACACA-----AATTAACCC[400]
Spar	TTCTTTCCGAGAAA---GCACGGCAACAAACA-----ACAGTCC-[400]
Smik	TTTCCCTCAAGGAATA-CGCACGGGTAACGAAA-----ACTGTCGG[400]
Sbay	TGGCTACCGGAAAAGACCCACAGCAGCTACTCATTGGCAACACCG[400]
Scer	-ATGCAACCCAGGAAAAAAATAGTCATATACGTAAGT-----[450]
Spar	-ATATGACCAAAAATAACCGCAATGGCAGTGAATGCAATTAT[450]
Smik	AATATGACCAAAAAATAATAAGAACG-TGGGAAACG-AAGTAT[450]
Sbay	CAAACGGCA-GTAAGAGGCCGAGTGAGGT-TTTTGTACGAAGGGAAAGCCC[450]
Scer	CGCTTTCATTGA [462]
Spar	CATTGATACGA [462]
Smik	CACTGATACGA [462]
Sbay	CAGCTTGGCCG [462]

Notion de protogènes

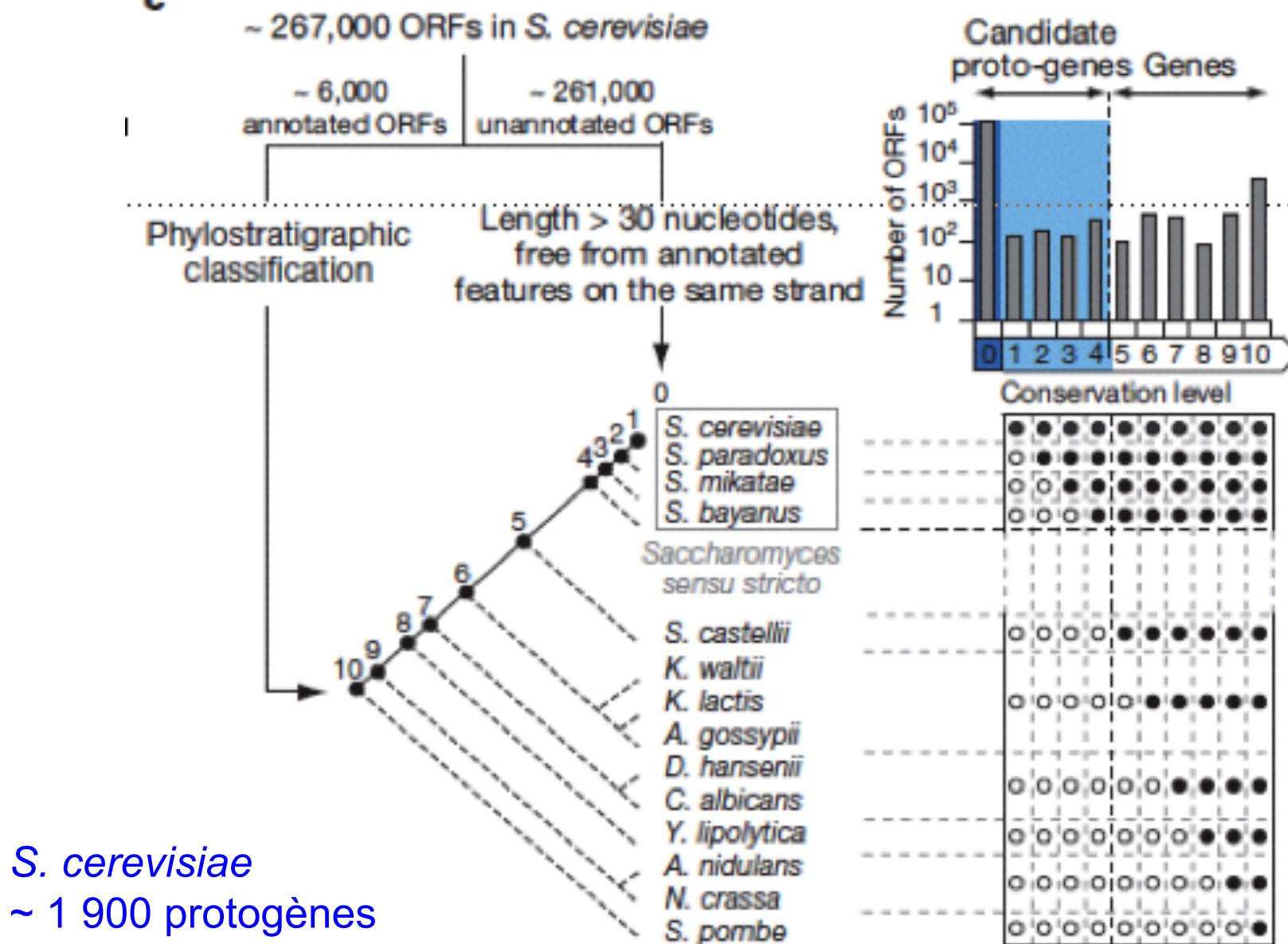
Carvunis et al. 2012 Nature 487: 370-374



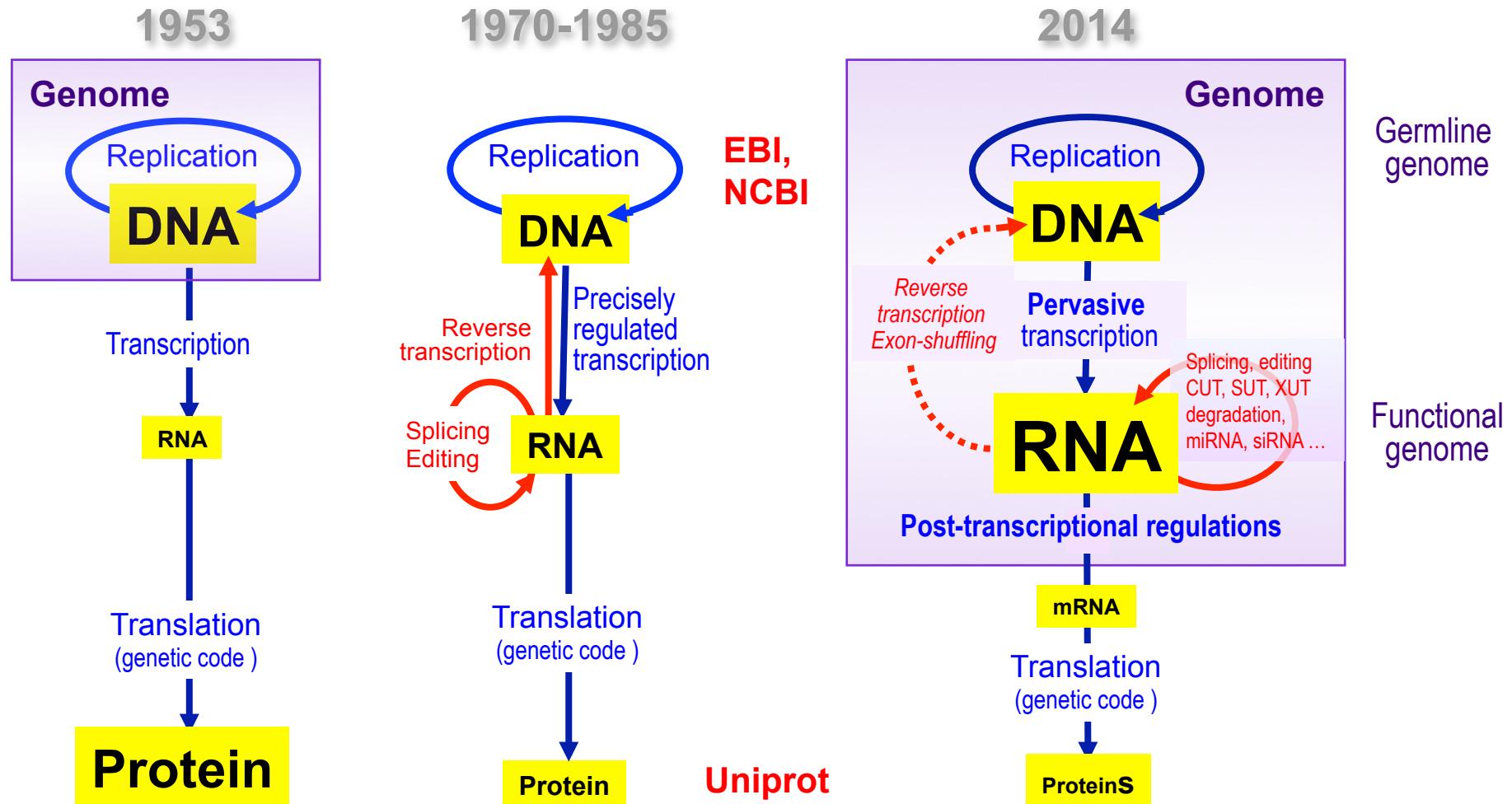
c

Notion de protogènes

Carvunis et al. 2012 Nature 487: 370-374



The central dogma of molecular biology, revisited

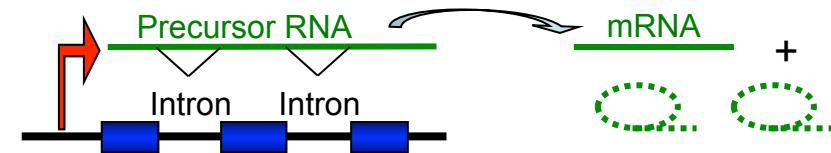


RNA surprises

1970 reverse transcription
Chicken virus



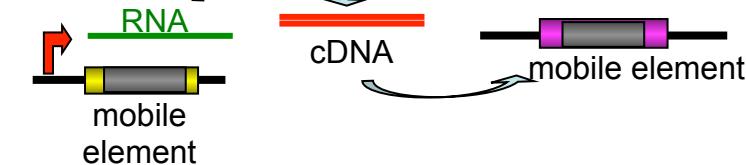
1977 introns
Mammalian virus, rabbit globin gene



1983 RNA catalysis
Tetrahymena nuclear intron, E. coli RNase P



1983 RNA editing
Trypanosoma

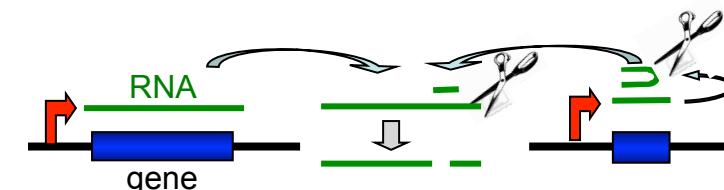


1985 retrotransposons
Saccharomyces cerevisiae



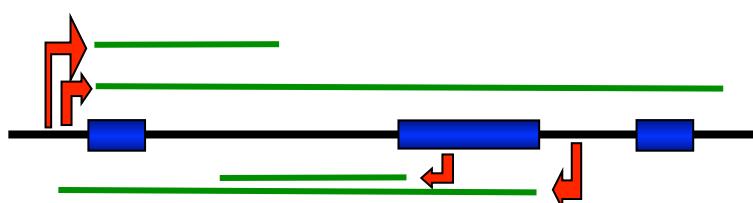
1993 cellular retrogenes
Drosophila melanogaster

2000 RNA interference
petunias, fungi, Caenorhabditis elegans



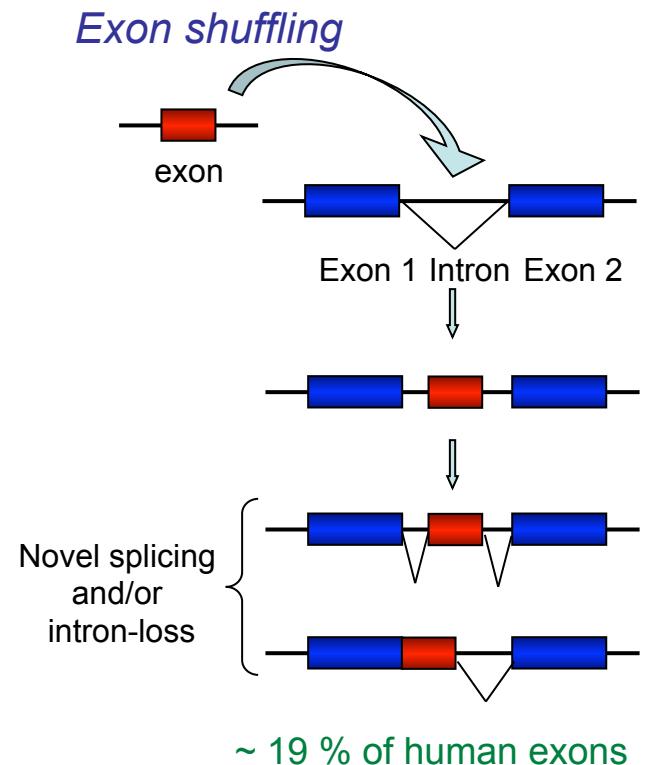
2002 micro RNAs
Caenorhabditis elegans

2012 pervasive transcription
S. cerevisiae, human





Walter Gilbert



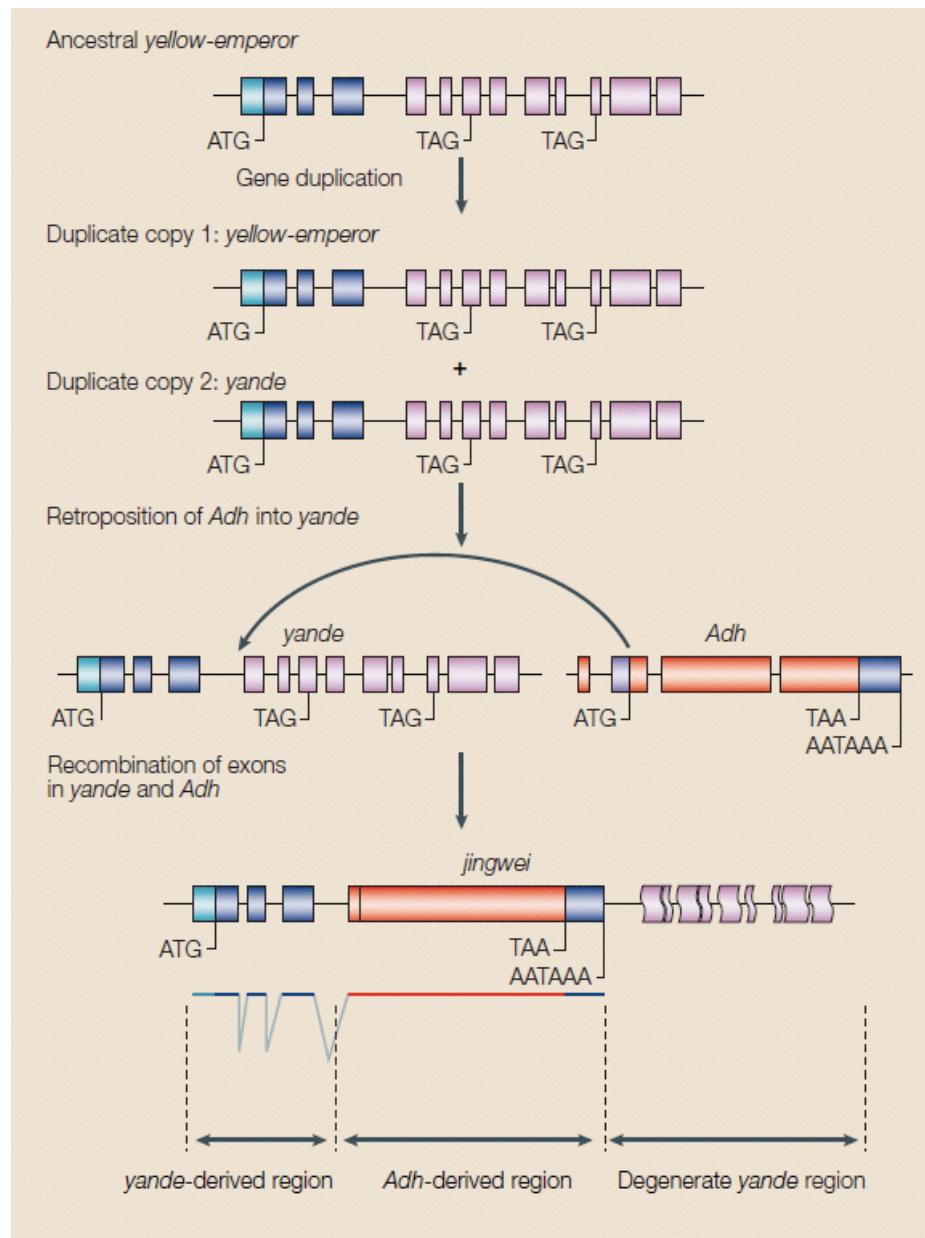
Retrogenes

≈ 1 % of
human exons

*Exonization of
mobile elements*

≈ 4 % of
human exons

RNA-mediated events



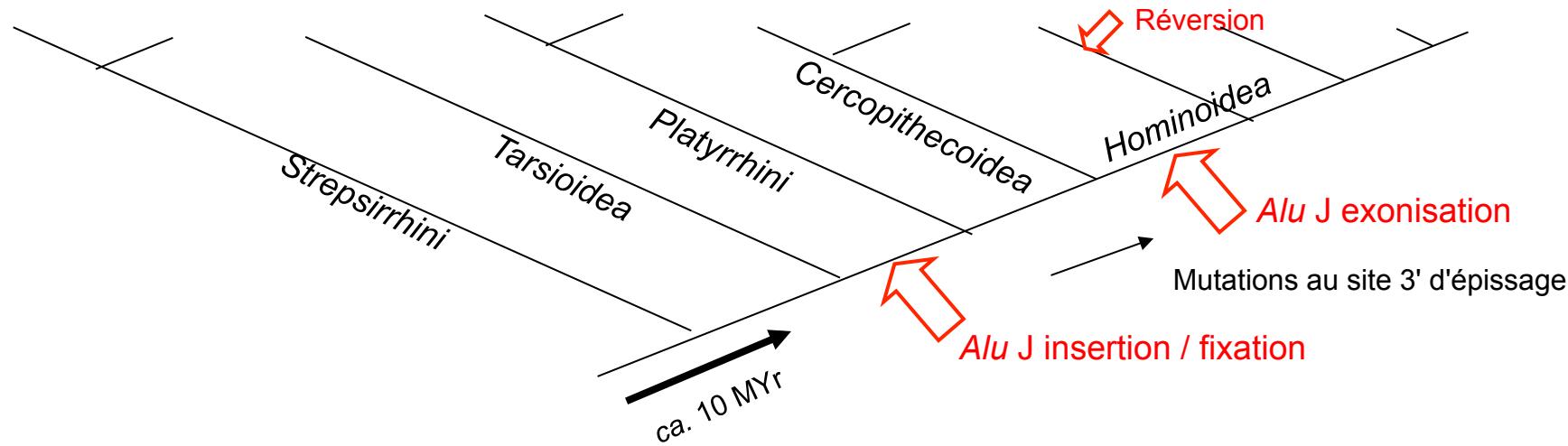
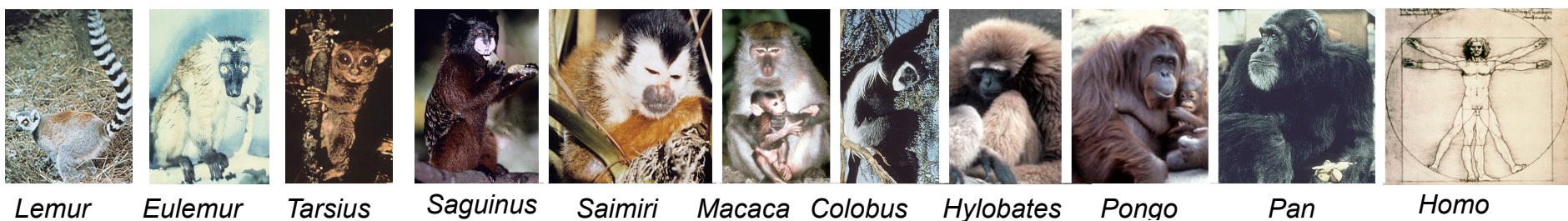
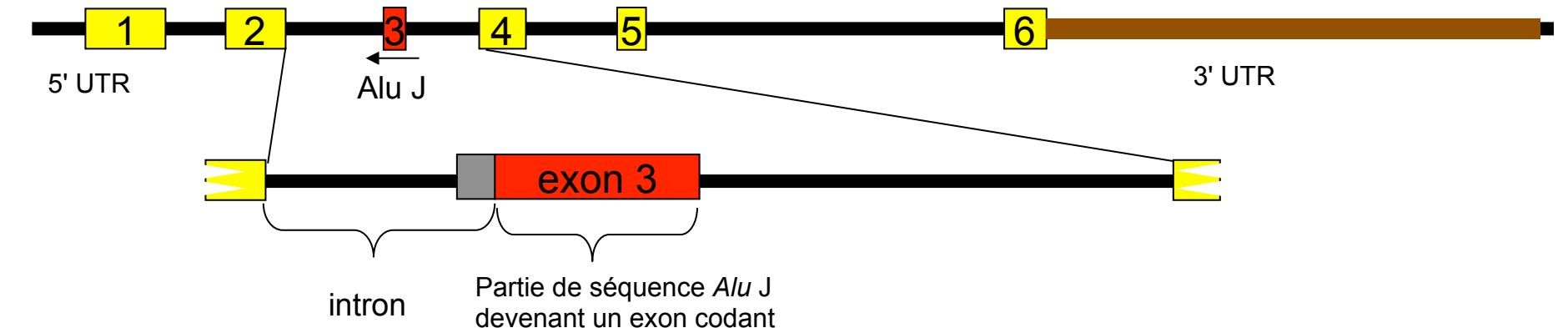
Long and Langley (1993) *Science* **260**: 91-95

Long et al. (2003) *Nature Genetics Reviews* **4**: 865- 875

Exonisation d'éléments mobiles

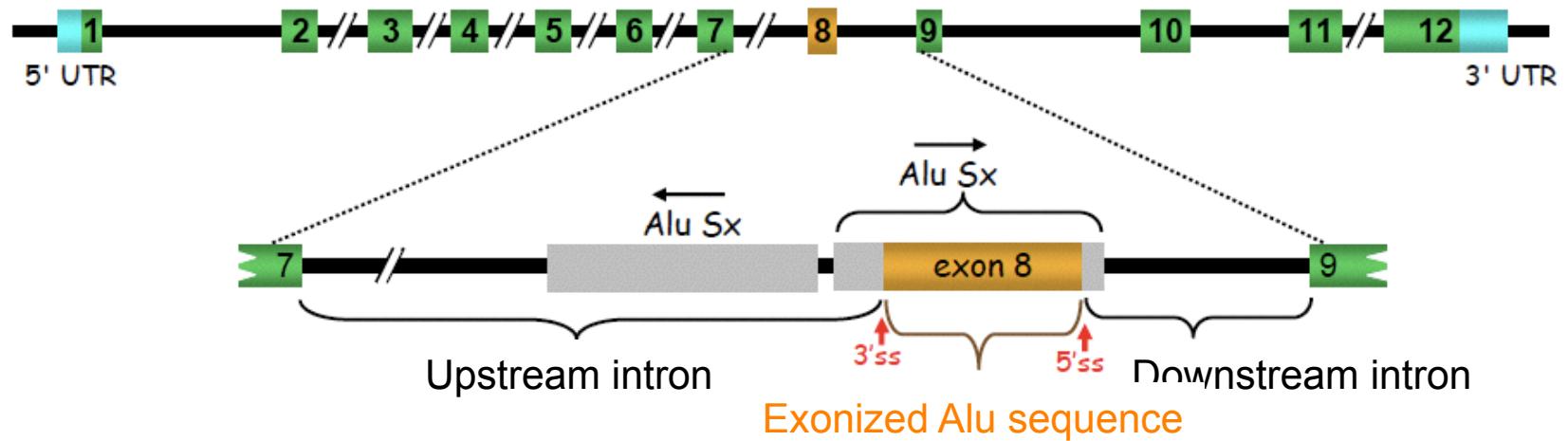
Gène humain *RPE2-1* Ribulose-5-phosphate-3-épimerase

Krull et al., (2005) Mol. Biol. Evol. 22, 1702-1711

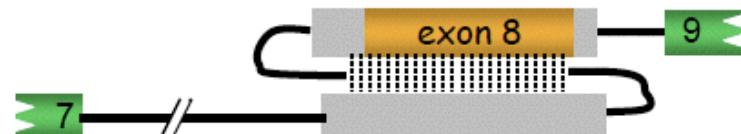


Exonization of mobile elements

Nuclear prelamin A Recognition Factor (Möller-Krull et al., 2008 J. Mol. Biol.)



A to I RNA editing generates the 3' SS and eliminates an in-frame stop codon



Exonization of mobile elements

