

CSE206B-L2

Genome rearrangements: WGD

Yeast Whole Genome Duplication?

Molecular evidence for an ancient duplication of the entire yeast genome

Kenneth H. Wolfe & Denis C. Shields

Department of Genetics, University of Dublin, Trinity College, Dublin 2, Ireland

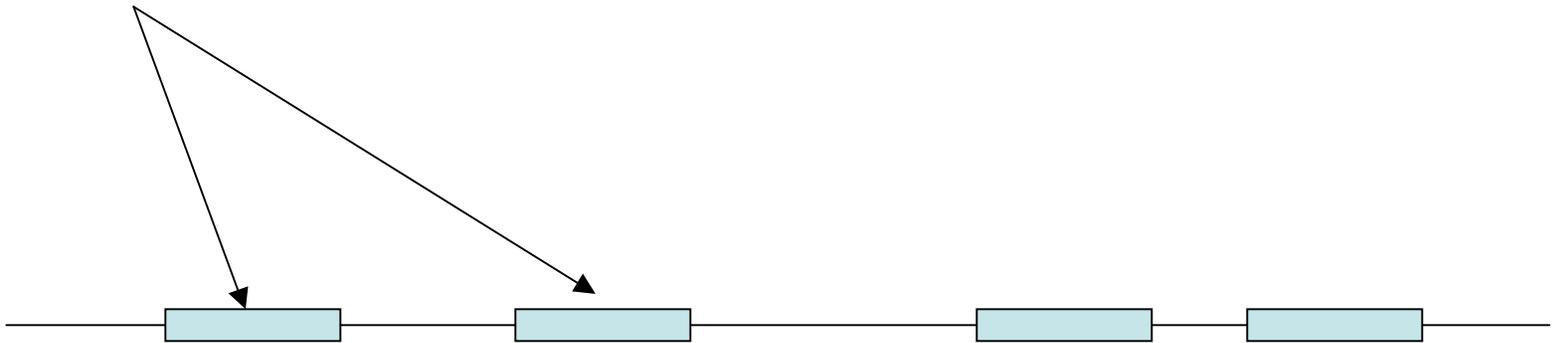
.....

be phenotypically redundant³⁻⁷. Here we show that the arrangement of duplicated genes in the *S. cerevisiae* genome is consistent with Ohno's hypothesis. We propose a model in which this species is a degenerate tetraploid resulting from a whole-genome duplication that occurred after the divergence of *Saccharomyces* from *Kluyveromyces*. Only a small fraction of the genes were subsequently retained in duplicate (most were deleted), and gene order

Prove or Disprove the WGD hypothesis

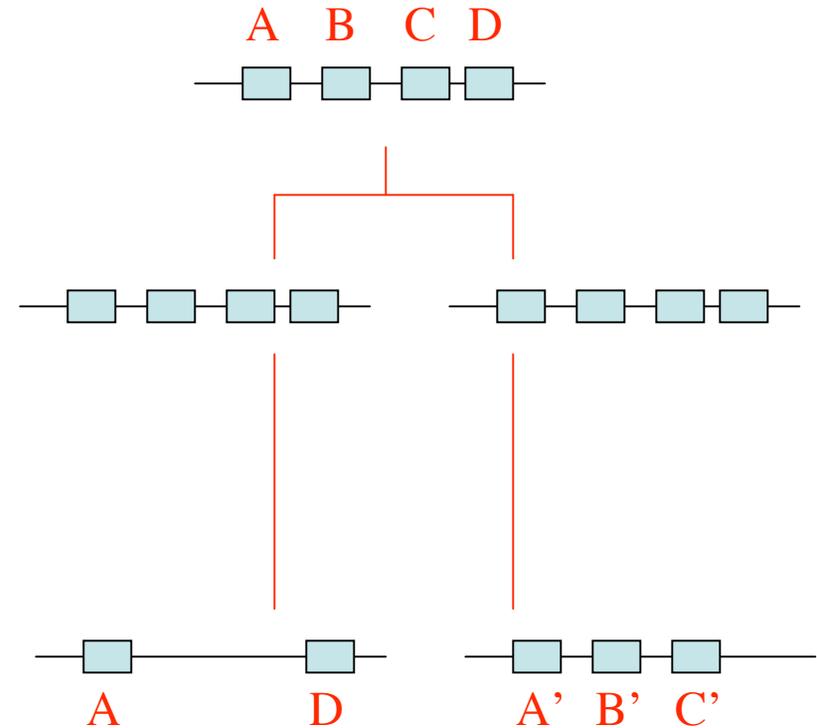
- First, some terms:

- DNA
- Proteins
- Genes



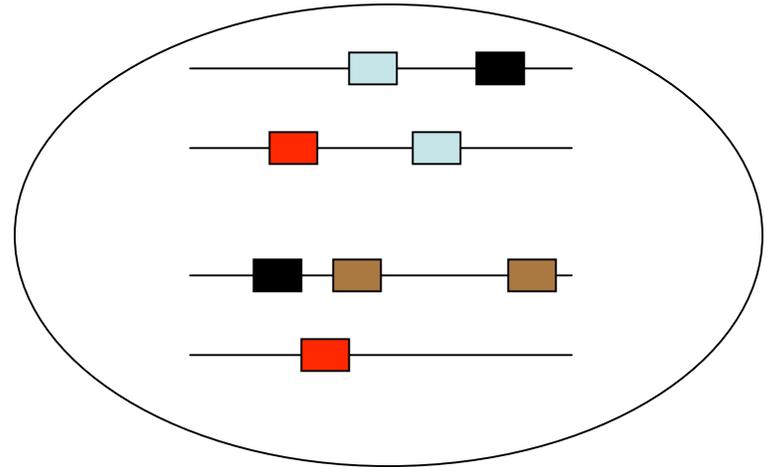
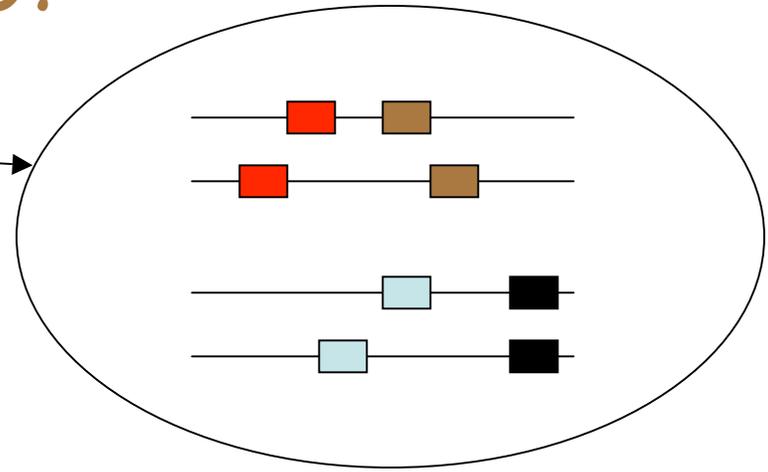
WGD

- Each ancestral chromosome split into two.
- Genes were rapidly deleted.
- A few duplicates remain



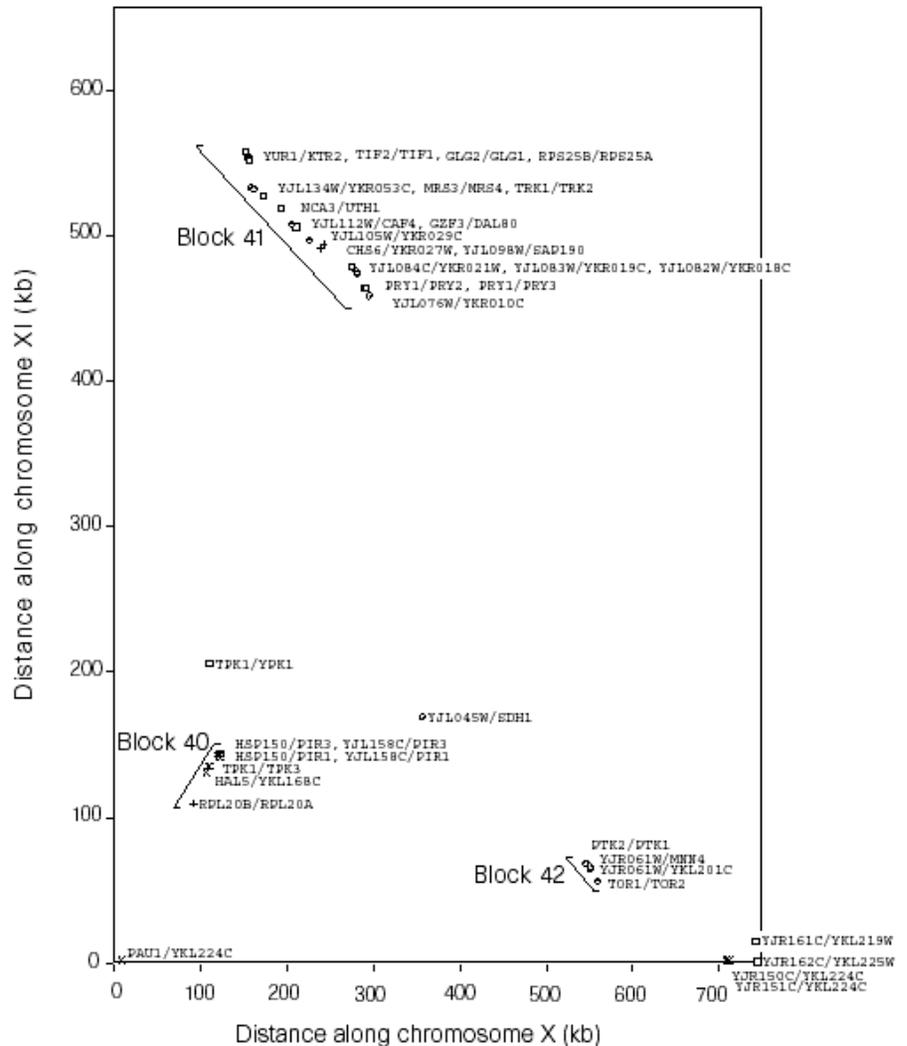
Are few duplicates enough to show WGD?

- Yes!
- Unfortunately, the genomes are also undergoing rearrangements.
- This complicates the picture



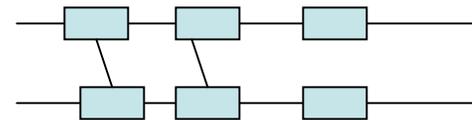
Evidence

- All *S. cerevisiae* proteins were compared against each other using BLASTP.
- 55 genomic regions were duplicated containing 376 duplicated genes



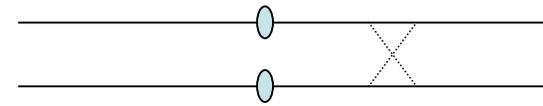
Duplicated regions

- What defines a duplicate region
 - BLAST score ≥ 200 for each pair
 - At least 3 pairs within 50kb each
 - Conservation of gene order and orientation (some local inversions allowed).
- 25% of the genes in duplicate regions were duplicated.
- 42% of all duplicate genes were in duplicated regions



How did the 55 regions arise?

- Independent duplications of each region.
- Alternatively, WGD followed by reciprocal translocations.
- 1. 50/55 regions have the same orientation w.r.t centromere.
- 2. Based on a Poisson distribution, 55 successive independent duplications would result in 7 triplicate regions. None are seen.
- Therefore, independent duplication is unlikely.



WGD or Not?

- The Wolfe paper concludes that yeast is an ancestral tetraploid. They conjecture that this happened when two diploid yeast cells fused to form a tetraploid.
- Sequence decay/deletions led back to diploidy.
- About 15% of the gens are still in duplicate form
- This conclusion was not widely accepted!

Genomic exploration of related gene families

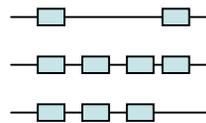
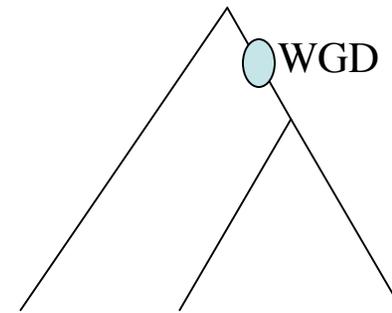
- Llorente et al., 2000 presented a different argument.
- Conducted survey sequencing of 13 yeast species.
- Looked at gene families, and found that the overall degree of gene redundancy seems conserved across all species.
- This includes species that predated the duplication, which should have a different distribution.
- Conclusion. No evidence for WGD.

Wong et al. 2002

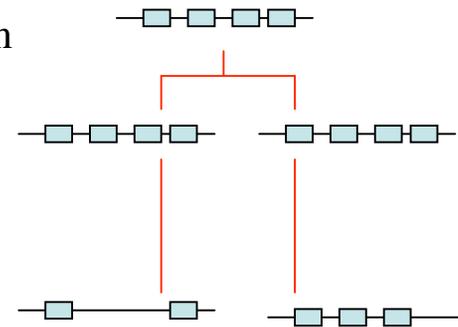
- They find even more compelling evidence of WGD from the 14 survey sequences!
- Can you suggest a comparative experiment?

Clues from pre-duplicated genomes

- Pairs of post-duplicated genomic regions should correspond to one pre-duplicated genomic region.



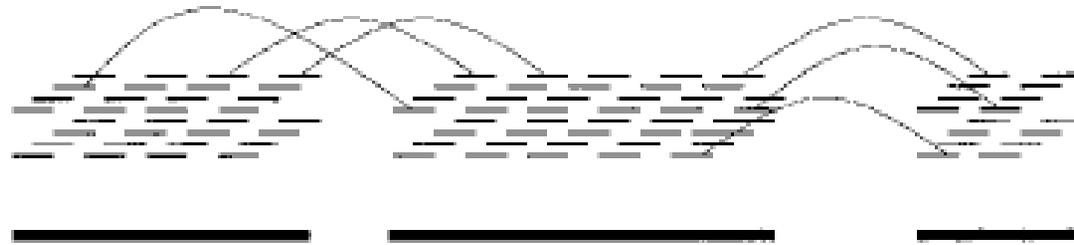
Pre-duplication



Post-duplication

- In reality, the picture is convoluted by local rearrangements. However, there is still considerable evidence remaining.

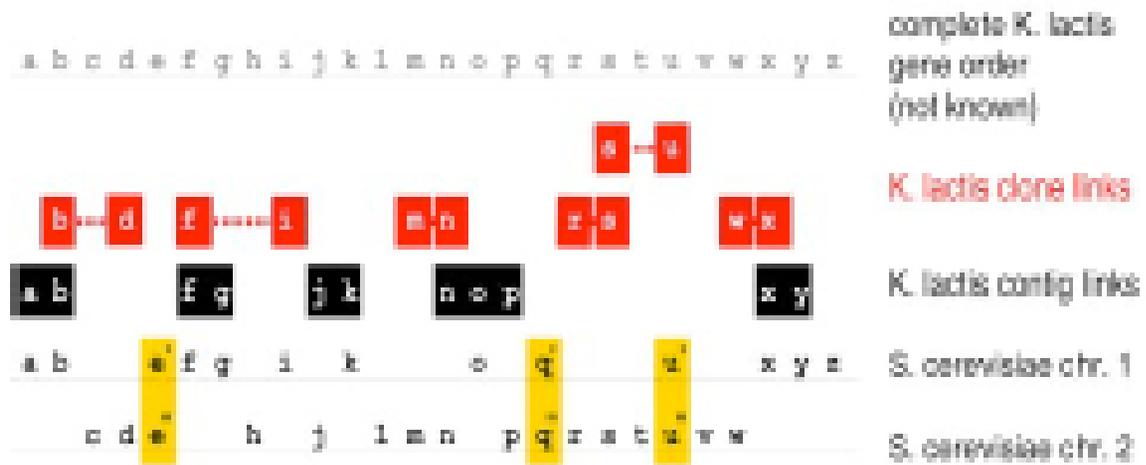
Genome sequencing & assembly



- Fragments are sequenced from the ends of large clones.
- Overlapping fragments are assembled (alignment and consensus generation) to form contigs.
- Clone mate-pairs are used to connect contigs.

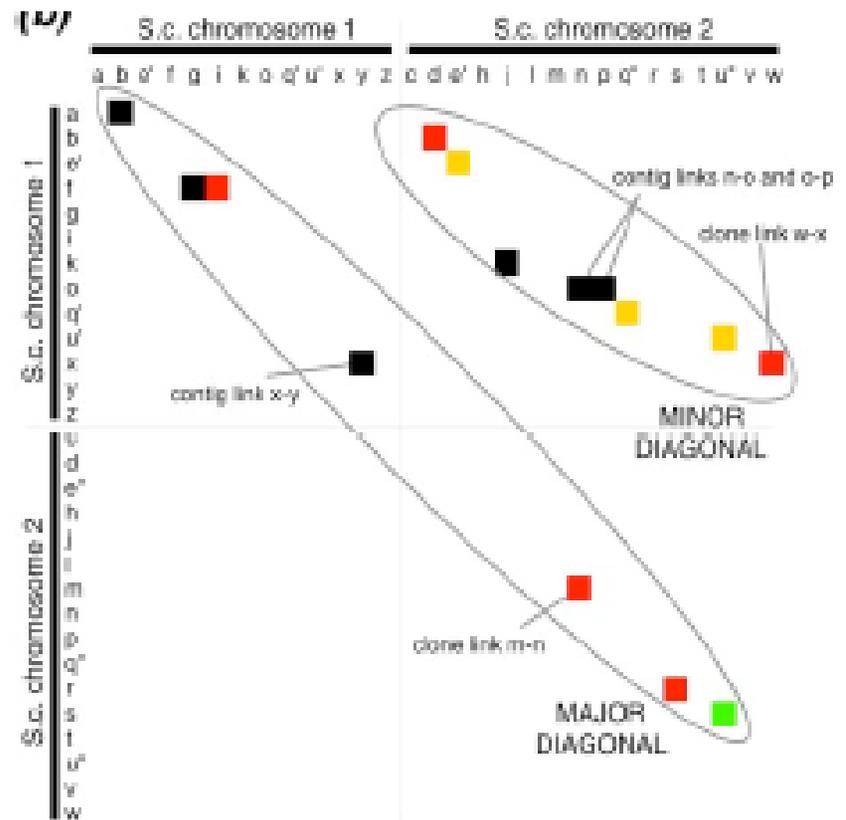
Pre-duplicated genome

(a)

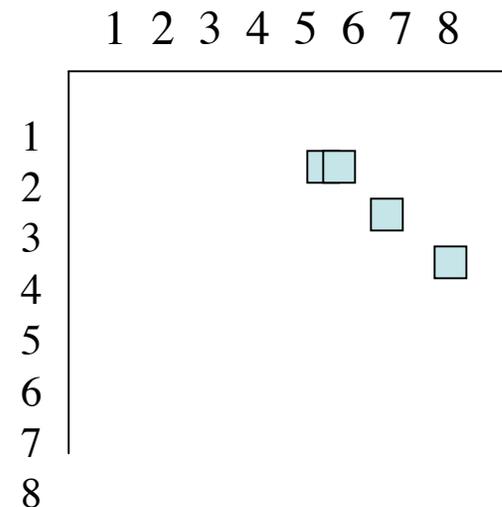
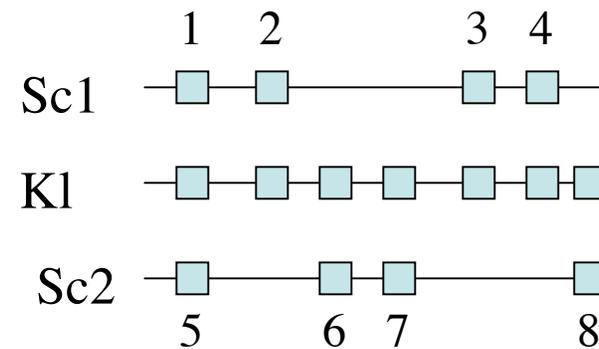


Comparing pre, and post WGD genomes

- Both axes represent *S. cerevisiae* genes
 - Black dots represent contig links in *K.lactis*
 - Red dot represent clone links
 - Yellow dots represent duplication in *S.cere.*
- What does this mean?



- Major diagonal represents synteny between Sc and Kl
- Minor diagonal band represents similarity between 2 Sc chromosomes inferred by Kl gene order



- 88 minor diagonals were identified.
- These diagonals covered 176 sections of *Sc* chromosomes
- 3900/5583 genes were included (70%)
- Almost no overlap between regions.
- Superimposing information about duplicated genes, the sister regions covered 82% of the genome.
- No overlap between sister regions rules out independent duplication

Final proof for WGD, 2004

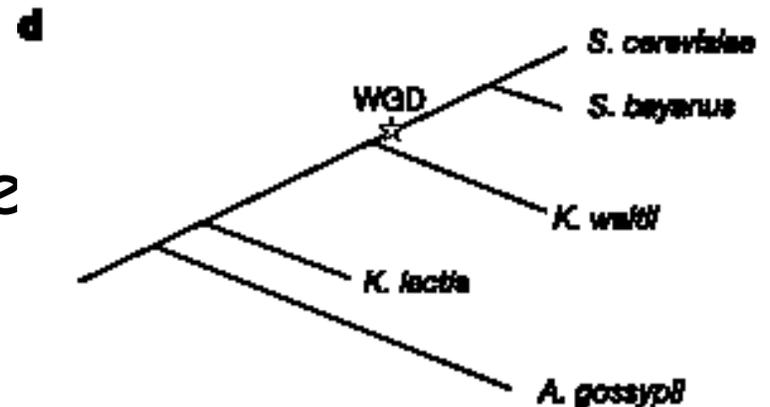
Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*

Manolis Kellis^{1,2}, Bruce W. Birren¹ & Eric S. Lander^{1,3}

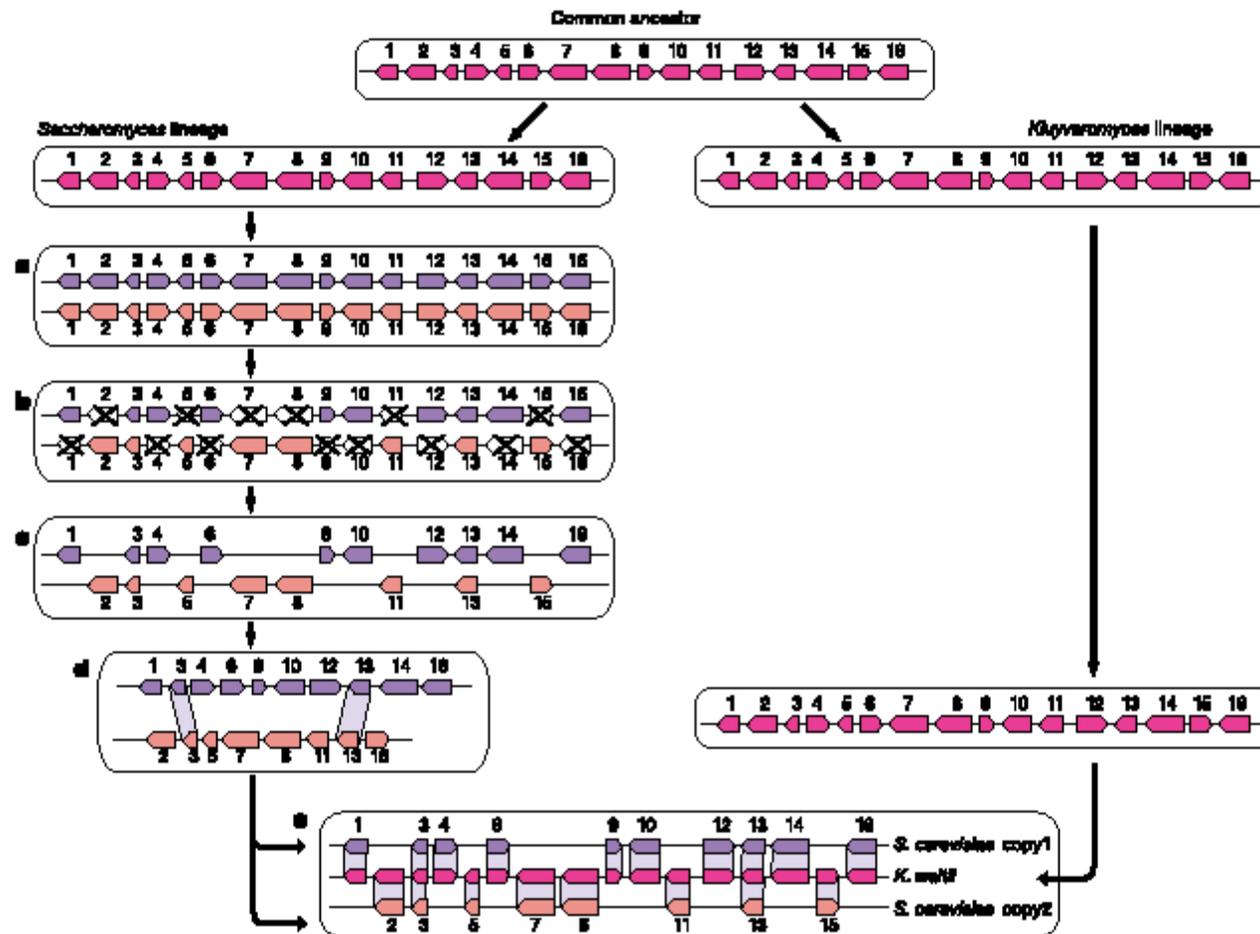
¹The Broad Institute, Massachusetts Institute of Technology and Harvard University, Cambridge, Massachusetts 02138, USA

²MIT Computer Science and Artificial Intelligence Laboratory, and ³Whitehead Institute for Biomedical Research, Cambridge, Massachusetts 02139, USA

- Complete sequence of a pre duplicated genome was generated.



Evidence is similar to Wong et al.



Doubly Conserved Synteny Blocks (DCS)

- 253 DCS blocks were identified containing 75% of *K. waltii* genes and 81% of *S. cerevisiae* genes
- A typical DCS block has 27 genes (largest block has 81 genes).
- DCS blocks are separated by ~3 genes on the average.
- ~1% of the *K. waltii* genome lies in segments that match 3 or more *Sc* regions. Mainly due to local rearrangements.
- In a DCS block 90% of *Kw* genes have a match in at least 1 of the 2 *Sc* regions.
- 47 Sister regions have no duplicated gene.

Analysis of duplicated genes

- Gene loss can occur by large segmental deletions, or individual gene deletions.
- DCS blocks indicate that gene loss occurred by many small deletions, and is balanced between the 2 sister regions (Average: 0.57-0.43)

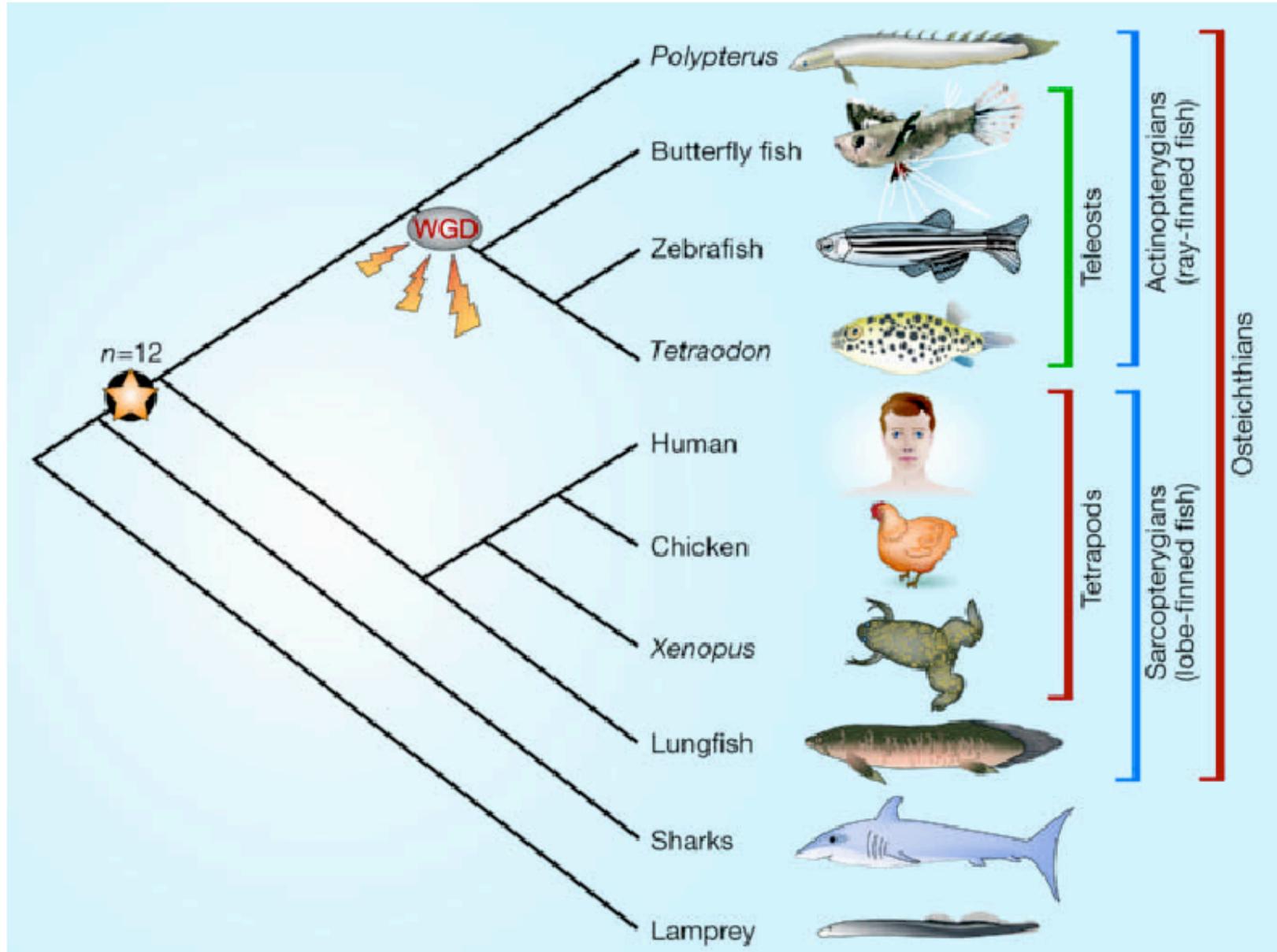
Accelerated divergence

- Ohno hypothesized that after duplication, one copy would preserve the original function, and the other copy would be free to diverge. Others argued that both copies would diverge.
- Kellis et al. Compared genes from *K. waltii* to corresponding pairs in *Sc*.
- 76 of 457 pairs show accelerated evolution. In 95% of the cases, acceleration was limited to one of the 2 paralogs.
- Deletion of the ancestral paralog is lethal in 18% of the cases.
- Deletion of a derived paralog is never lethal.

WGD in Pufferfish: 2004

Genome duplication in the teleost fish *Tetraodon nigroviridis* reveals the early vertebrate proto-karyotype

Olivier Jaillon¹, Jean-Marc Aury¹, Frédéric Brunet², Jean-Louis Petit¹, Nicole Stange-Thomann³, Evan Mauceli³, Laurence Bouneau¹, Cécile Fischer¹, Catherine Ozouf-Costaz⁴, Alain Bernot¹, Sophie Nicaud¹, David Jaffe³, Sheila Fisher³, Georges Lutfalla⁵, Carole Dossat¹, Béatrice Segurens¹, Corinne Dasilva¹, Marcel Salanoubat¹, Michael Levy¹, Nathalie Boudet¹, Sergi Castellano⁶, Véronique Anthouard¹, Claire Jubin¹, Vanina Castelli¹, Michael Katinka¹, Benoît Vacherie¹, Christian Biéumont⁷, Zineb Skalli¹, Laurence Cattolico¹, Julie Poulain¹, Véronique de Berardinis¹, Corinne Cruaud¹, Simone Duprat¹, Philippe Brottier¹, Jean-Pierre Coutanceau⁴, Jérôme Guzy⁸, Genis Parra⁶, Guillaume Lardier¹, Charles Chapple⁶, Kevin J. McKernan⁹, Paul McEwan⁹, Stephanie Bosak⁹, Manolis Kellis³, Jean-Nicolas Volff¹⁰, Roderic Guigó⁶, Michael C. Zody³, Jill Mesirov³, Kerstin Lindblad-Toh³, Bruce Birren³, Chad Nusbaum³, Daniel Kahn⁸, Marc Robinson-Rechavi², Vincent Laudet², Vincent Schachter¹, Francis Quétier¹, William Saurin¹, Claude Scarpelli¹, Patrick Wincker¹, Eric S. Lander^{3,11}, Jean Weissenbach¹ & Hugues Roest Crollius^{1,4}



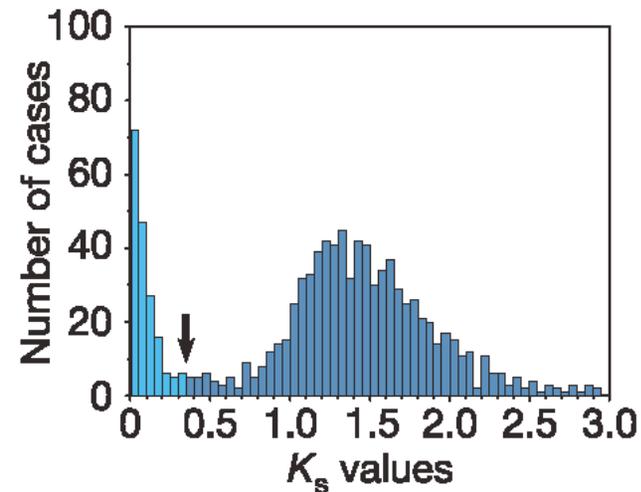
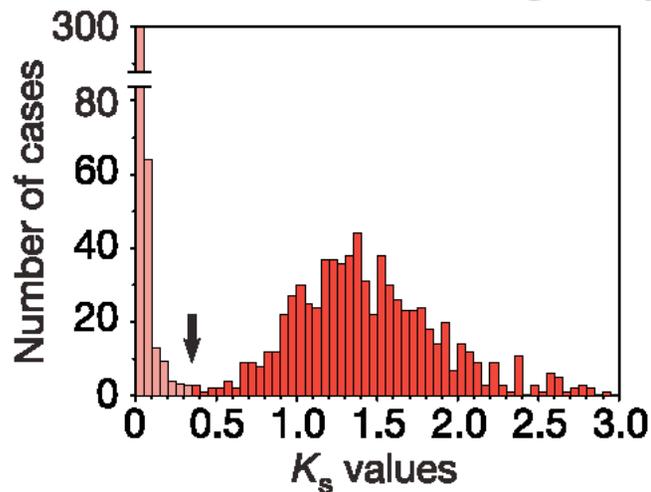
2 Main arguments for WGD

- Evidence for a number of 'paired' duplicate regions. The regions do not overlap and must not show higher redundancy
- Evidence for DCS when compared against a pre-duplicated 'nearby' genome.

Selecting ancient duplications

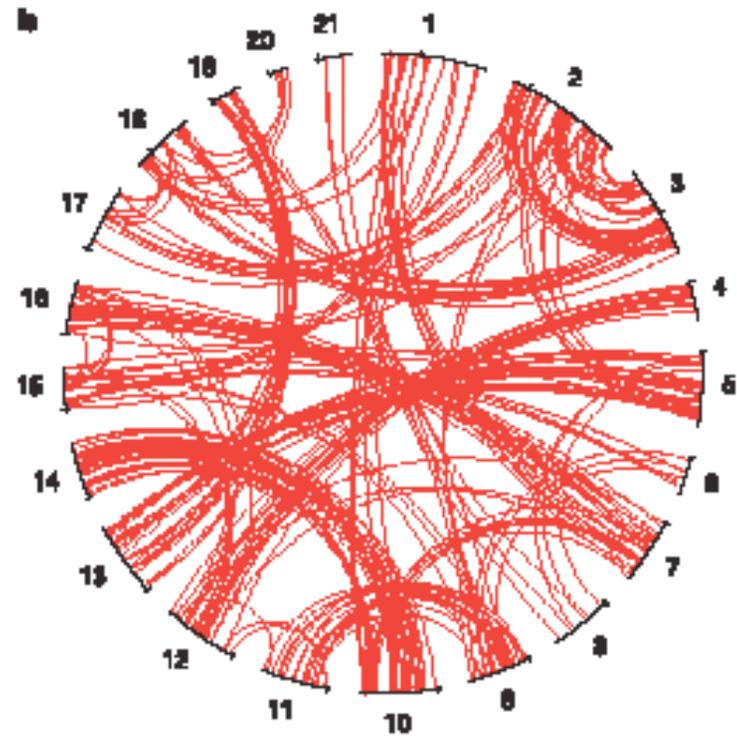
Duplicated genes

- 1078 pairs in tetra, and 995 pairs in fugu
- based on K_s , 75% are ancient duplications before tetra-fugu separation

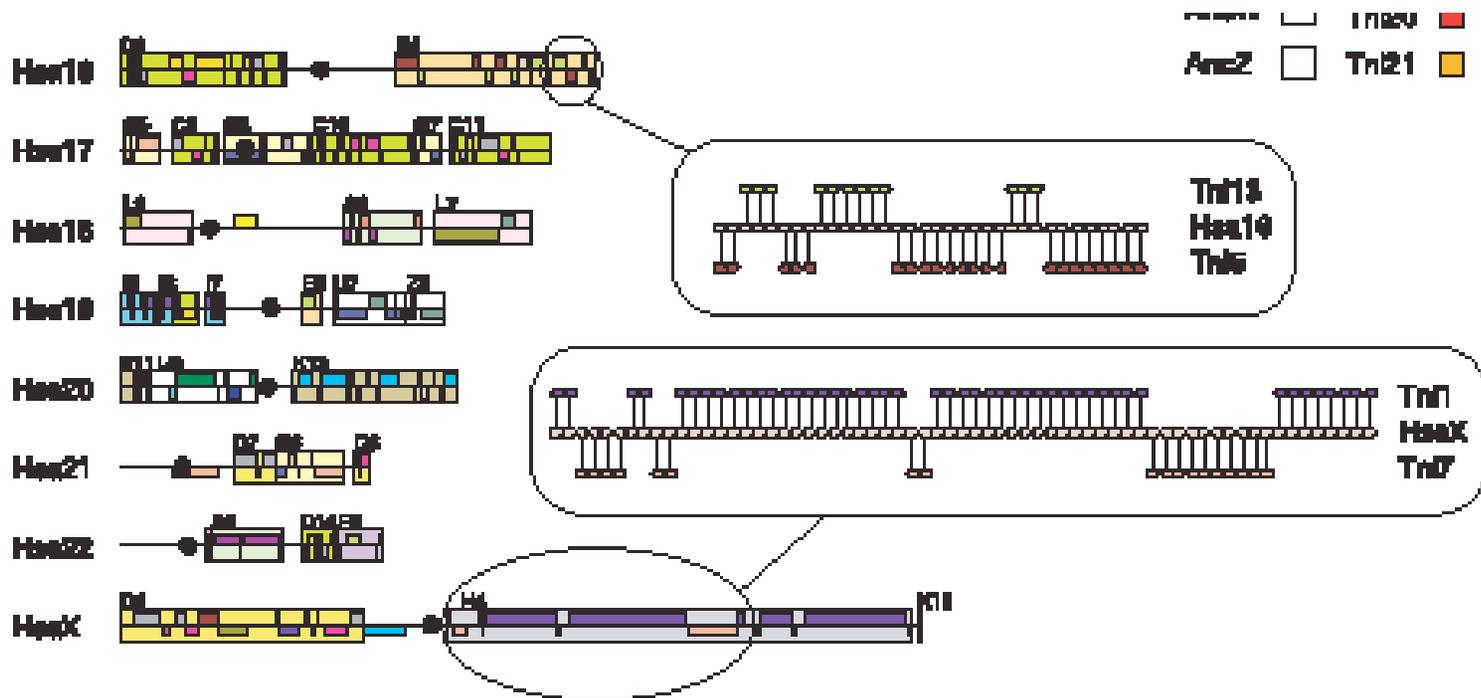


Evidence 1

- Duplicated sections of the chromosome.



Evidence 2

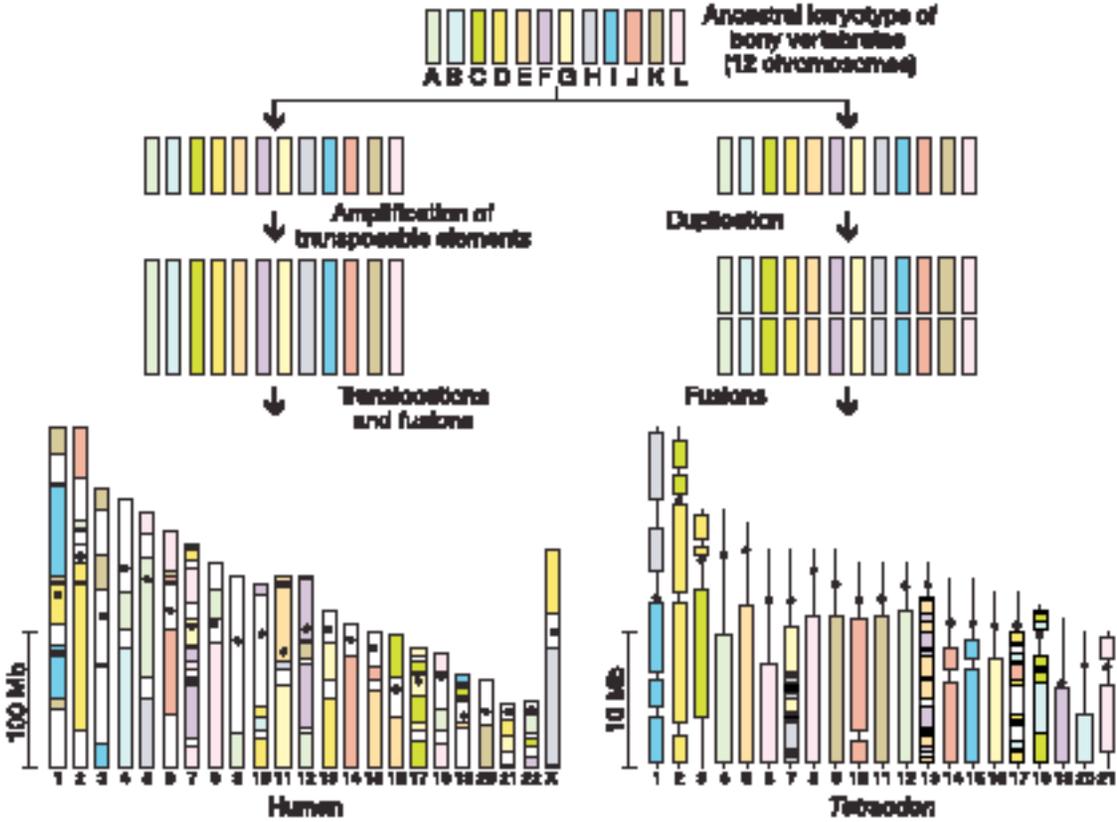


- Compared against human and mouse

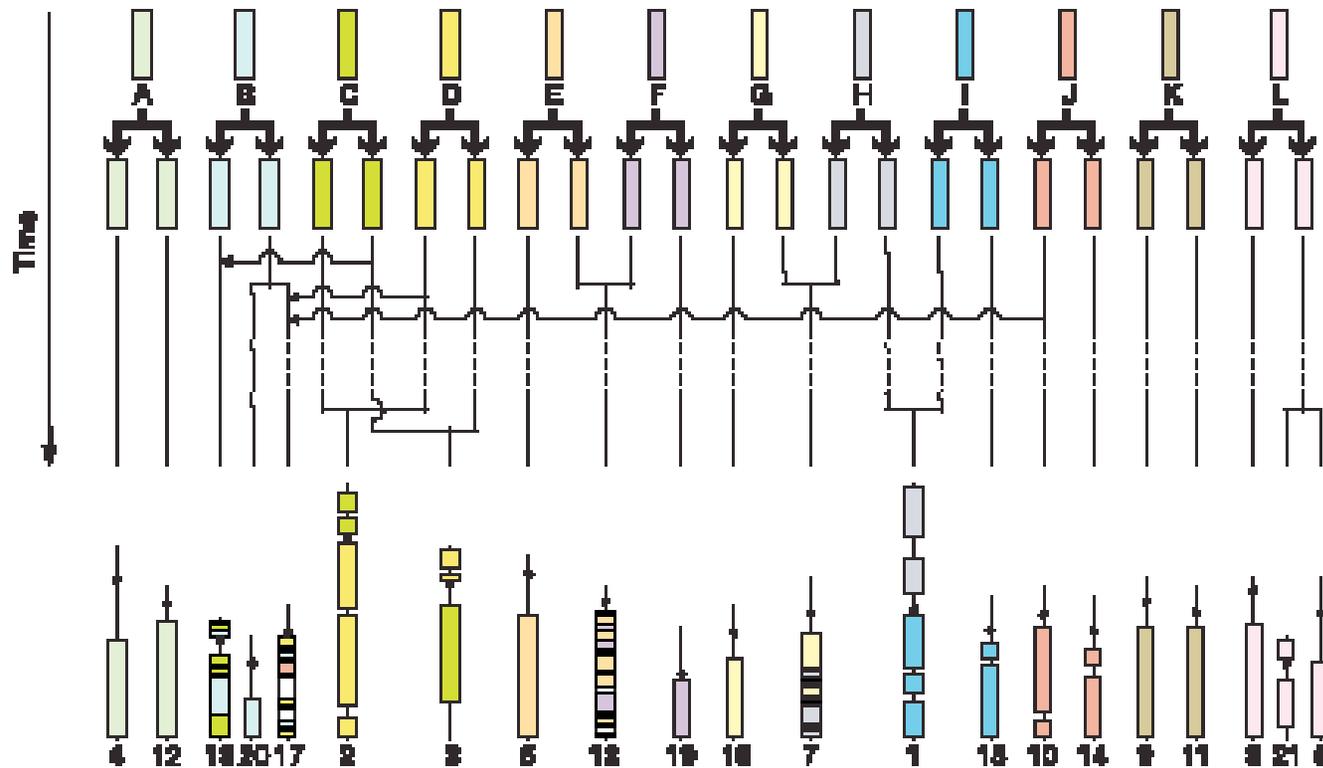
Composition of the ancestral genome?

- What did the pre-duplicated ancestor look like?
- Can you postulate the architecture of the current (human and tetraodon genomes) genomes in terms of the common ancestor?
- What was the sequence of events?

Model for distribution of ancestral segments in current Hs and Tn genomes



An evolutionary scenario for Tn



Genome rearrangement scenarios

- In subsequent classes, we will cover some of the theoretical approaches to studying these genome rearrangements.

Other topics

How many genes do we have?

- There appear to be about 30,000–40,000 protein-coding genes in the human genome—only about twice as many as in worm or fly. However, the genes are more complex, with more alternative splicing generating a larger number of protein products.

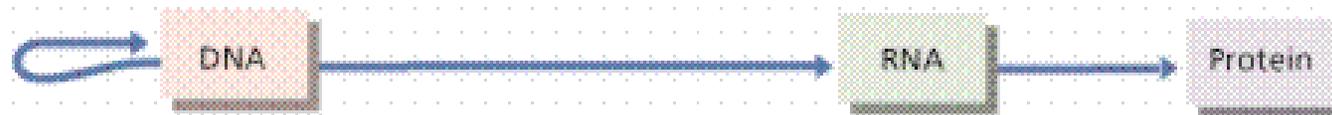
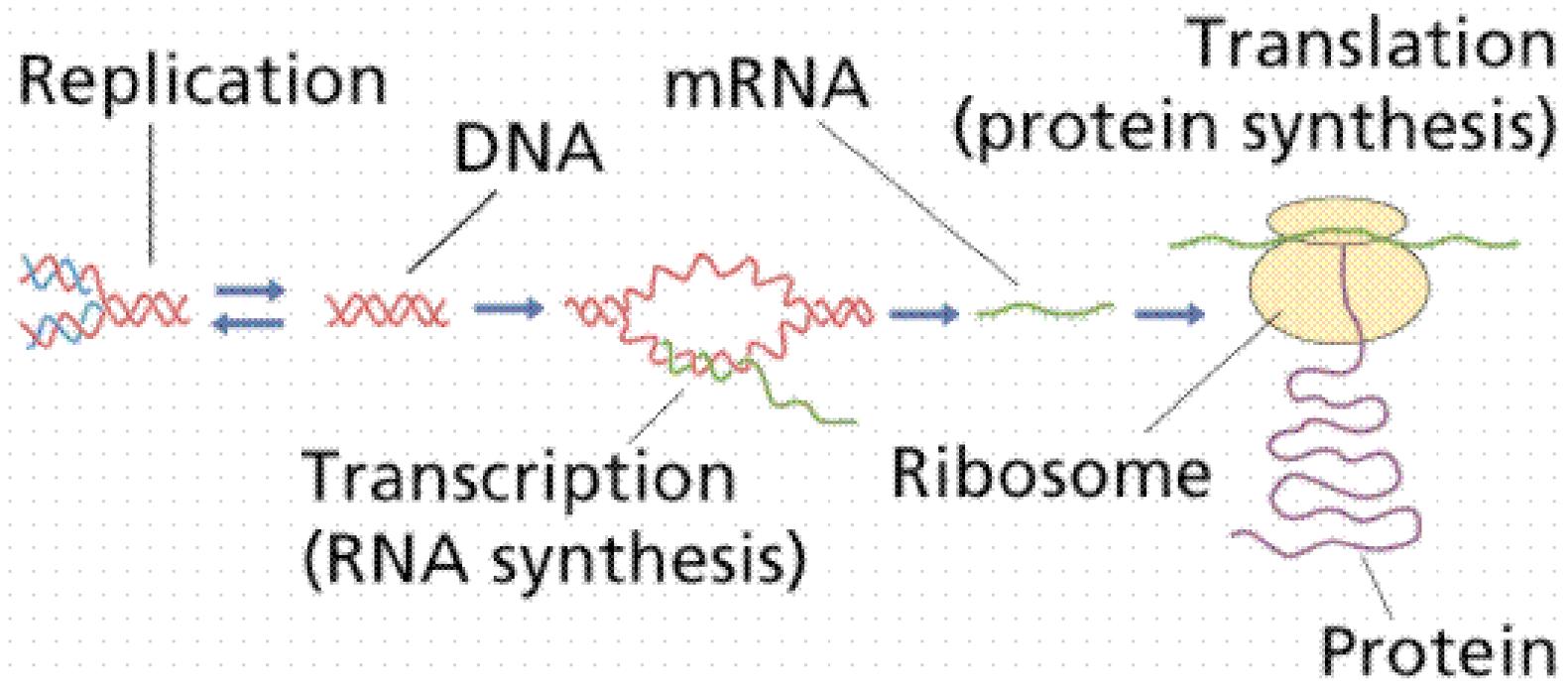
Nature 2001

Science

pressed. Conservative criteria, requiring at least two lines of evidence, were used to define a set of 26,383 genes with good confidence that were used for more detailed analysis presented in the subsequent sections.

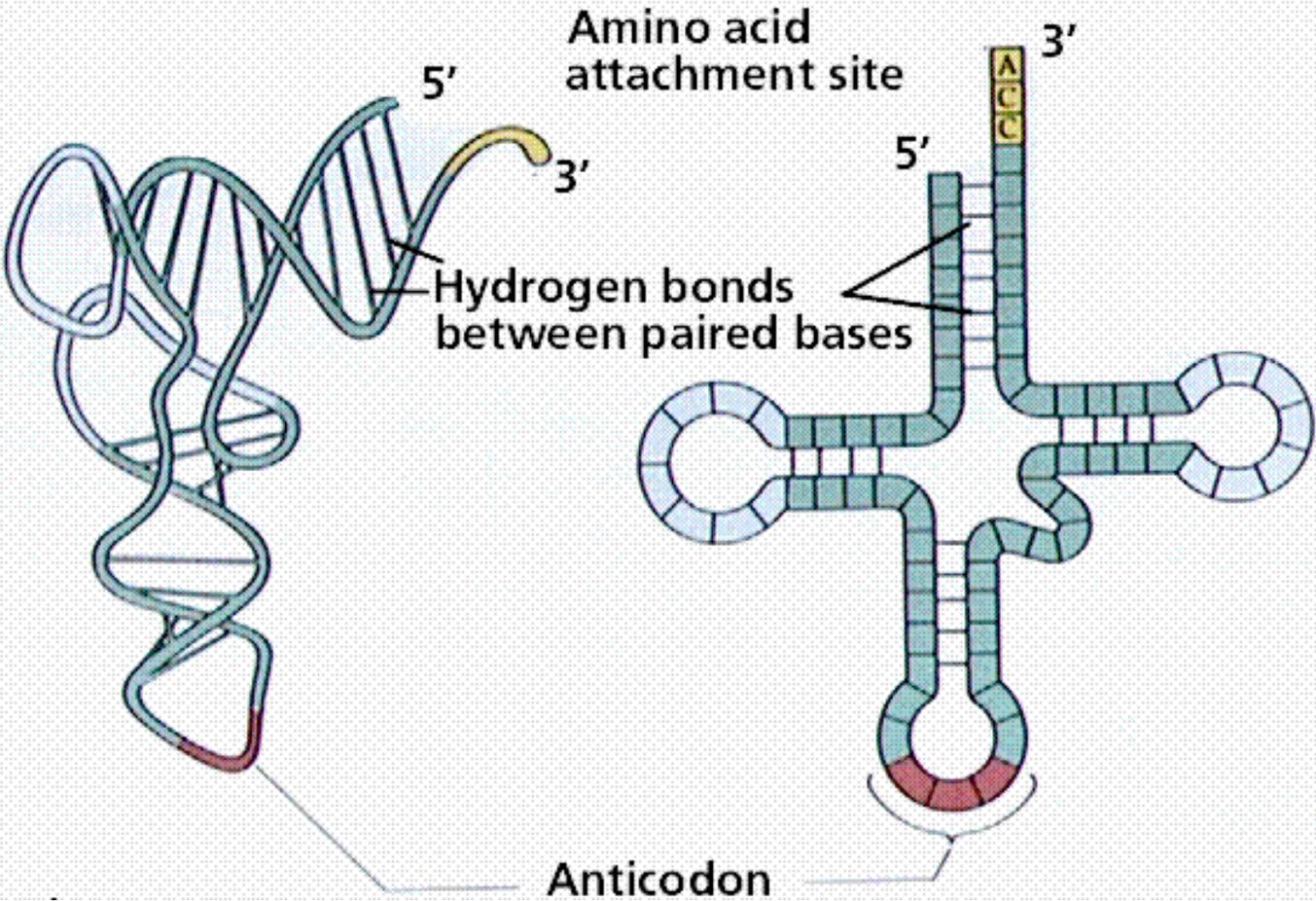
The Modern RNA World hypothesis

- The estimate of the number of genes is based largely on *computational* experiments. These methods works well for *large, highly expressed, protein coding genes*.
- “Could it be that a large class of genes has gone undetected because they do not make proteins?..”
 - Eddy: *Nature Reviews Genetics* 2001
- If true, what function can these genes have?



Roles of ncRNA

- mRNA is translated. All untranslated RNA are lumped as ncRNAs
- However, the Ribosome (the translating machinery) is composed largely of RNA (rRNAs).
- tRNAs encode the template for translation.
- Splicing machinery in the nucleus (spliceosome) contains snRNA.
- Many other nc RNAs (snoRNA, siRNA, miRNA) are being found on a regular basis. Some of these act as anti-sense regulators of gene expression (RNAi).
- The sequence specific binding makes RNA a more efficient regulator than proteins.

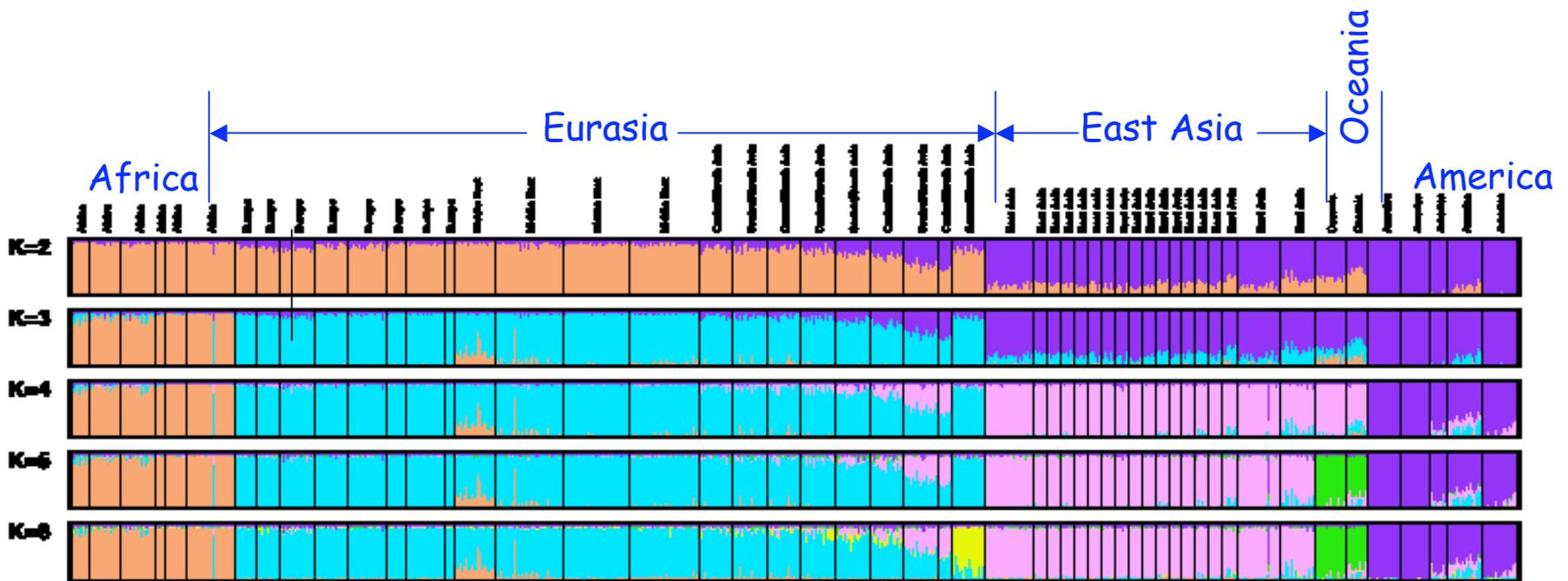


Computational ncRNA genefinding

- Understanding & Computing RNA structure
- De novo: Predict signals for transcription initiation, termination, processing
 - Eukaryotic ncRNAs are transcribed by different polymerases.
 - Approach feasible in microbes
- Comparative genome analysis
 - Compensatory mutations preserve structure
 - Predict structure from alignment
 - Searching with homologs

Population Structure

- 377 locations (loci) were sampled in 1000 people from 52 populations.
- 6 genetic clusters were obtained, which corresponded to 5 geographic regions (Rosenberg et al. Science 2003)



Population Genetics

- What is it about our genetic makeup that makes us measurably different?
- These genetic differences are correlated with phenotypic differences
- With cost reduction in sequencing and genotyping technologies, we will know the sequence for entire populations of individuals.
- Here, we will study the basics of this polymorphism data, and tools that are being developed to analyze it.

Repeats

- Roughly 45-50% of the human genome contains transposable repeats.
- Like rearrangements, these are fossil records, and can offer unique insights into evolution.
- We will discuss tools for identifying, classifying, and analyzing repeats.