Gene and genome duplication

CUSO/SIB Course Next Generation Comparative Phylogenomics 2014

Marc Robinson-Rechavi

http://bioinfo.unil.ch/ marc.robinson-rechavi@unil.ch @marc_rr

L | Université de Lausanne Département d'écologie et évolution



Mechanisms of duplication:

•Polyploidy

whole genome duplication

- •Aneuploidy
- •Duplicative transposition
- Local tandem duplication

small-scale duplication, or gene duplication

SIB

Département d'écologie

et évolution

Swiss Institute of Bioinformatics



Detecting duplications

11 ... : 0

UNIL Université de Lausanne Département d'écologie et évolution



Discuss with your neighbor

• What is more frequent: point mutations or duplications?

Mmil

UNIL | Université de Lausanne Département d'écologie et évolution



Duplication may be more frequent than single nucleotide subsitution in the human genome



Can we trust Ks histograms?



Département d'écologie et évolution





Unil

UNIL | Université de Lausanne Département d'écologie et évolution





SSD age distributions are characterized by a saturation peak



Duplication great and small



UNIL Université de Lausanne Département d'écologie et évolution



Duplication can lead to lineage specific expansion of gene families





Figure 1 Phylogenetic analysis of selected eukaryotic lineage-specific expansions. Groups supported by a bootstrap value >70% are colored pink for Drosophila melanogaster, red for Homo sapiens, orange for Caenorhabditis elegans, green for Arabidopsis thaliana, and yellow for Schizosaccharamyces pombe. (A) Prolyl hydroxylases. (B) Small molecule kinases (Ch stands for choline kinase). (C) Patched-like protein. (D) MAP-Kinases. (E) P450 family hydroxylases. (F) MBDAT membrane acyltransferases. (At) Arabidopsis thaliana; (Bs) Bacillus sublik; (Ce) Caenorhabditis elegans; (Dd) Dictyostelium discoideum; (Dm) (Drosophila melanogaster; (Hs) Homo sapiens; (Pbcv1) Paramecium bursaria Chlorelia virus 1; (Rs) Rabtonia solanacearum; (Sa) Staphylococcus aureus; (Sc) Saccharomyces cerevisiae; (Sm) Sinorhizobium meliotė; (Sp) Schizosaccharomyces pombe. Complete tree descriptions (full lists of GI numbers or gene names, and bootstrap values) are available in the Supplementary Material online at ftp://ncbi.nlm.nih.gov/ pub/aravind/expansions, and http://www.genome.org..

Unil

Département d'écologie et évolution



Whole genome duplications



UNIL Université de Lausanne Département d'écologie et évolution





All pairs of proteins with BLASTP scores greater than or equal to 200 are plotted at the position of their genes on the two chromosomes. Ty elements have been omitted.

Symbols indicate gene orientations: +, W (Watson strand; left-to-right transcription) on both chromosomes; times, C (Crick strand; right-to-left transcription) on both; squares, C on chromosome X but W on chromosome XI; circles, W on chromosome X but C on chromosome XI.

© 2009 SIB

UNIL Université de Lausanne Département d'écologie et évolution

11 ...; 0





 χ^2 relative to random distribution of transcriptional orientation and gene order P < 0.0001 less triplicates than expected by chance (Poisson) P = 0.001

≈ 85% loss of duplicatesafter whole genomeduplication

55 duplicate regions were identified containing 376 pairs of homologous genes. Amino-acid sequence identity between the pairs: 24% to 100%, with a mean of 63%. Criteria used to define a duplicate region:

- (1) BLASTP high scores of greater than or equal to 200 for each gene pair (these have an associated significance of $P = 10^{-18}$ or less)
- (2) \geq 3 pairs of homologues with intergenic distances \leq 50 kb on each chromosome
- (3) conservation of gene order and orientation (with allowance for small inversions within some blocks).

Duplicated regions on average 55 kb long, 6.9 duplicate gene pairs -> span 50% of the genome.

Département d'écologie et évolution





		Bootstrap		90% Confidence	Amino-acid	Branch length		ath	
Block	Gene 1 / Gene 2	support (%)	Age	intervial	sites	A	В	c	Outgroup
3	NTH2/NTH1	100	0.70	(0.59 - 0.82)	714	0.052	0.131	0.115	Candida albicans
3	GAL1/GAL3	99	0.55	(0.41 - 0.76)	468	0.100	0,115	0.130	Homo sapiens
37	YCK1/YCK2	97	0.48		343	0.036	0.046	0.024	Sch. pombe
36	HXK1/HXK2	94	0.77	(0.62 - 0.95)	484	0.037	0.147	0.107	Sch. pombe
45	EXG1/SPR1	90	0.80	(0.67 - 0.97)	418	0.051	0.166	0.270	Candida albicans
10	CRY1/CRY2	89	0.40		137	0.015	0.014	0.007	Horno sapiens
21	SIR2/HST1	45	0.86	(0.66 - 1.00)	402	0.028	0.167	0.182	Caenorhabdilis elegans
28	CYC7/CYC1	34	0.82		106	0.017	0,110	0.057	Sch. pombe
30	YGR043C / TAL1	24	1.00		332	0.000	0.231	0.149	Homo sapiens
Genes with apparent gene conversion									
8	RPS10A / RPS10B	100	0.04		233	0.042	0.000	0.004	Sch. pombe
20	SSB1 / SSB2	100	0.06	(0.01 - 0.11)	613	0.049	0.002	0.004	Candida albicans
37	RPL41B / RPL41A	100	0.00		105	0.054	0.000	0.000	Candida tropicalis

Duplications which can be dated occured after the S. cerevisiae / Kluyveromyces speciation



UNIL | Université de Lausanne Département d'écologie et évolution

Unil



A few take home points:

- the ancient whole genome duplication was unsuspected before genome sequencing + bioinformatics
- synteny information provided key evidence
- most genes were lost after the whole genome duplication





Of the importance of comparing genomes

UNIL | Université de Lausanne Département d'écologie

et évolution



Model of WGD followed by massive gene loss predicts gene interleaving in sister regions



UNIL | Université de Lausanne Département d'écologie et évolution

11 ... : 0





Gene and region correspondence with *K. waltii* reveals Whole Genome Duplication

a- Each region of *K. waltii* (coloured by number of matches: white, 0; grey, 1; black,
2; yellow, > = 3) shows conserved gene order with two regions of S. cerevisiae (coloured by chromosome number). Spacing between *S. cerevisiae* genes is set to match *K. waltii* chromosomal positions. Vertical blue bars denote centromeres.

b- Doubly conserved synteny region showing duplicate mapping of centromeres (black circles). All sixteen *S. cerevisiae* centromeres show such duplicate mapping with *K. waltii* centromeres. This DCS region also illustrates that we can reliably recognize anciently duplicated segments even in the absence of any remaining two-copy genes. Evidence of WGD comes from gene interleaving and 2:1 mapping with orthologous *K. waltii* segments. The segments containing intervening genes are deleted, resulting in condensed sister regions. Kwal, *K. waltii*; Scer, *S. cerevisiae*.

NIL | Université de Lausanne Département d'écologie et évolution





Duplicated blocks in S. cerevisiae

Kellis et al (2004) Nature 428: 617-24

UNIL | Université de Lausanne Département d'écologie et évolution



A few take home points:

• data alone is not sufficient to prove whole genome duplication, good bioinformatics is needed

no obvious link between genome duplication and organism complexity
duplication can be proven even after secondary loss (e.g. there is only one gene left)



JNIL | Université de bausanne Département d'écologie et évolution



Also in plants

Unil

UNIL | Université de Lausanne Département d'écologie et évolution



Arabidopsis genome:

DNA based: All five chromosomes aligned to each other in both orientations using MUMmer results filtered to identify all segments at least 1,000 bp in length with at least 50% identity. -> 24 large duplicated segments of 100 kb or larger, comprising 65.6 Mb or 58% of the genome.

Protein based: TBLASTX to identify collinear clusters of genes residing in large duplicated chromosomal segments.

-> 67.9 Mb, 60% of the genome.

sequence conservation of the duplicated genes varies greatly:

6,303 (37%) of the 17,193 genes in the segments classified as highly conserved (E < 10 -30) 1,705 (10%) with less significant similarity up to E < 10-5.



Genomic map of duplicated blocks in Arabidopsis

The two copies of each putative duplicated block (e.g., 1a and 1b) are shown. Color denotes age class (red, A; blue, B; green, C; purple, D; orange, E; gray, F). Centromeres are shown with black circles, and ribosomal DNA with white circles. Direction of arrowhead indicates the predominant relative orientation of duplicated cORFs within each block (right, direct; left, inverted). Landmarks are given at 200 composite ORF intervals.



DNIL | Université de Lausanne Département d'écologie et évolution

11 0



Table 1. Features of the five age classes of duplicated blocks. d, is the minimum change in amino acids per between dispersed duplicated cORFs, averaged among all cORFs. Retained duplicates, ratio of presently duplicated to inferred ancestral cORFs. Block size, mean number of cORFs (including singletons) per copy.

Age class	No. of blocks	d _A	Retained duplicates	Block size	Estimated age (Mya)
A*	2	0†	0.90	12.5	0
В	1	0.21†	0.38	14.5	50
С	35	0.451	0.15	149.5	100
D	36	0.571	0.13	128.0	140
E	23	0.711	0.11	60.0	170
F	6	0.84†	0.09	57.8	200

*Probable artifact of genome assembly.
†Median among all duplicated genes.
‡Mean of best-fit normal
distribution to block medians.

Overlaps between the blocks

duplicates of different "ages"

-> many block duplications in the past evolution of Arabidopsis?

11 ... : 0



Sensitive methods are important



UNIL Université de Lausanne Département d'écologie et évolution



© 2009 SIB



Sib Swiss Institute of Bioinformatics

304 non hidden duplications with AHDoRetotal of 3571 anchor points+ 1607 tandem repeats, involving 4193 individual genes.



Overview of the chromosomal location of all multiplicons detected in the Arabidopsis genome

Example of a multiplicon in which nonhidden duplications can be observed between all three segments involved



NIL | Université de Lausanne Département d'écologie et évolution





Conclusion: 3 whole genome duplications in Arabidopsis evolution

NIL | Université de Lausanne Département d'écologie et évolution

11,



i-ADHoRe

Input: genes from one or several genomes -> matrix of homology by BLAST

clusters defined as significant by probability of finding another anchor by chance, given the total density of homologous pairs

iterative procedure to find multiple duplications (multiplicons):

search for homology to previously defined blocks of high score



i-ADHoRe



UNIL | Université de Lausanne Département d'écologie et évolution

Unil



Example of an Arabidopsis multiplicon

Segments shown clockwise in the order in which they were added to the multiplicon:

Segments 1 and 2 were detected first as a pair of collinear segments that were consequently aligned to create a profile.

This profile was then used to detect segment 3, which, in turn, was aligned and added to the



© 2009 SIB



Département d'écologie

et évolution

SIB

Based on the similarities with any of the other individual segments, the homology with segment 8 is far from statistically significant because there are too few anchor points and the number of intervening genes is too high. However, if we consider the seven other segments in the multiplicon (of which the mutual homology has already been established) together as a profile, we see that segment 8 shares in total eight genes with the profile.









© 2009 SIB

UNIL | Université de Lausanne Département d'écologie et évolution

Unil

Gene order alignment of collinear regions conserved over a large phylogenetic distance (human–chicken).



NIL | Université de Lausanne Département d'écologie et évolution



Discuss with your neighbor

- What is more frequent: point mutations or duplications?
- Where in animal evolution were there genome duplications?



Département d'écologie et évolution



History of an idea:

Susumo Ohno (1970): observation of higher haploid DNA content in mammals than Ciona

-> hypothesis of two whole genome duplications which would account for "higher complexity"



http://www.mun.ca/biology/desmid/brian/BIOL3530/DB_Ch15/BIOL2900_EvoDevo.html

11...:0

Département d'écologie et évolution


Two rounds of whole genome duplication in vertebrates





Garcia-Fernandez & Holland (1994): One Hox cluster in Amphioxus -> duplications from 1 to 4 were vertebrate specific



Holland (2003) J Struct Funct Genomics 3: 75-84

Spring (1997): many examples of gene families with 1 ortholog in Drosophila, 4 in human but

anecdotal evidence

no phylogenetic or syntenic test

11 ...:0



And more genes in fishes

Unil



More Hox in fishes





Unil





UNIL | Université de Lausanne Département d'écologie et évolution

Unil



To understand vertebrates, study amphioxus

JNIL Université de Lausani

et évolution

Département d'écologie











It's not because it's amphioxus that it's always simple

> UNIL | Université de Lausanne Département d'écologie

> > et évolution



Expansion of innate immunity genes in amphioxus (right) compared with the human system (left).



Bref, whole genome duplications abund

Unil





Kellis et al 2004 Nature 428:617-624



Jaillon et al 2006 Nature 444, 171-178



Jaillon et al. 2004 Nature 431:946-957



Jaillon et al. 2007 Nature 449, 463-467



Yu et al 2005 PLoS Biol 3:e38



Putnam et al 2008 Nature 453:1064-1071

Unil UNIL | Université de Lausanne Département d'écologie et évolution







et évolution



Let's work!

Unil



pick one of FRE1-8; examine region in: http://ygob.ucd.ie/

pick a human gene from list, examine region in: http://www.genomicus.biologie.ens.fr/

Unil



Consequences of duplication

11 ...;



Discuss with your neighbor

- What is more frequent: point mutations or duplications?
- Where in animal evolution were there genome duplications?
- Two reasons why we should distinguish orthologs and paralogs





Ohno (1970): after duplication, either one copy gains a new function, or one copy is lost



Lynch lab (1999-2001): probability of new functional mutation before loss too low relative to frequency of duplicate genes



Département d'écologie et évolution



© 2009 SIB

DDC model = Duplication Degeneration Complementation

Original model based on regulatory regions (i.e. gene expression patterns):



© 2009 SIB Force et al. (1999) Genetics 151: 1531-45

UNIL | Université de Lausanne Département d'écologie et évolution

11 ... : 0









UNIL | Université de Lausanne Département d'écologie et évolution

Unil





Nature Reviews Genetics



Forms of functional divergence





Discuss with your neighbor

- What is more frequent: point mutations or duplications?
- Where in animal evolution were there genome duplications?
- Two reasons why we should distinguish orthologs and paralogs
- What is gene function?

11 ... : 0



Sequence divergence: Biochemical function

	35		H			ю	H
Hardin, Harrish Hardin, Harrish Hardin, Harrish Hardin, Harrish Hardin, Harrish Hardin, Harrish Hardin, Harrish Harrish Harrish			T. BOARD				
	1.4		-	1		1111	
NAME AND A DESCRIPTION OF THE PARTY OF THE P		1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1		1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1			
	**************************************	1 11					

© 2009 SIB

Escriva et al. 2006 PLoS Genet 2: e102

Département d'écologie et évolution





Transcriptional Activity and Binding Selectivity of Vertebrate RARs

Unil





Staining of embryos indicates expression of mRAR α (A), mRAR β (B), and mRAR γ (C) in mouse embryos at E9; of xRARα (G), xRARβ (H), and xRARy (I) in stage 30 Xenopus embryos, and of AmphiRAR (M) in 20 h old amphioxus larvae. Schematic representations are shown of the expression territories of mRARs (D–F), xRARs (J–L), and AmphiRAR (N) in mouse, Xenopus, and amphioxus embryos, respectively. Regions with high levels of expression are red and those with lower levels of expression are pink. Arrowheads indicate regions in mouse and Xenopus embryos where the RAR expression cannot be correlated with AmphiRAR expression and can be described as "new expression territories."



11 ... 0







Not all genes are kept in duplicate





Are genes kept in duplicate biased relative to selective pressure and evolutionary rates?

genes are identified.

(C. elegans, S cerevisiae)



Ka = original rate of gene + possible change due to duplication



Outgroup species in which representative pairs of orthologs are identified.

Ka Drosophila / Anopheles = independent of changes after duplication in yeast or nematode







Comparison	S. cerevisiae		C. elegans	
	Spearman Correlation (r _{ab})	Partial Correlation Coefficient (r _{abc})	Spearman Correlation (r _{ab})	Partial Correlation Coefficient (r _{ale})
Class ^a versus CAI Class ^a versus K _S	0.40***	0.37*** -0.07m	0.09* 0.08 ^{ns}	0.12*
K _s versus CAI	-0.31***	-0.26***	-0.27***	-0.28***

Significance was tested for the direct and partial correlation coefficients using the statistics $t_s = r\sqrt{\frac{n-2}{1-s^2}}$ and $t_s = r\sqrt{\frac{n-2}{1-s^2}}$, respectively, where n is the sample size, m is the number of variables held constant, and r is the rank correlation coefficient (Sokal and Rohlf 1995).

"For this parameter, representative pairs were given a value of either 0 (for a singleton) or 1 (for a duplicate).

NS, nonsignificant; *, p = 0.05; **, p = 0.01; ***, p = 0.001.

DOI: 10.1371/journal.pbio.0020055.t002

1	Inil	



In yeast, biases differ between WGD and smaller scale duplications (SSD):

Table 1. Evolutionary	rate estimates for	the three sets o	f study genes	obtained by	measuring t	the divergence I	between outgroup
orthologs							

Gene set	Nemat	ode outgroup lineages ^b	Fly outgroup lineages ⁶		
	K _A (all genes)	K _A (without ribsomal genes)	K _A (all genes)	K _A (without ribsomal genes)	
Singletons	0.095	0.097	0.075	0.077	
WGD duplicates	0.061 ***	0.089 *	0.045 **	0.061 *	
SSD duplicates	0.079 **	0.087 *	0.063 **	0.070 *	

* The gene studied included singltons, WGD duplicates and SSD duplicates. Both types of duplicates appear to arise from conserved genes.

*Pairs of outgroup orthologs were identified in C. elegans and C. briggsae where possible (see methods in supplementary material online) and mean divergence between these pairs for each group, both with and without ribosomal proteins is shown.

"Similar mean divergence estimates for representative pairs from D. melanogaster and D. pseudoobscura are shown.

^dSignificance levels are * P<0.05, ** P<0.01 for Mann-Whitney U-test, compared with the set of singleton genes.

GO function		Including ribosomal	genes	Without ribosomal genes			
	Singletons	WGD pairs	SSD pairs	Singletons	WGD pairs	SSD pairs	
Catalytic	48.54%	45.21%	48.74%	51.64%	63.06% *b	54.19% ns	
Binding	24.73%	21.00%	26.38%	26.19%	21.66% ns	28.77% ns	
Transcription regulator	10.48%	9.13%	4.77%	11.24%	12.74% ns	5.31% *	
Structural molecule	10.39%	30.59%	11.81%	9.54%	4.46% *	10.89% ns	
Transporter	8.90%	3.20%	9.81%	3.92%	3.18% ns	1.96% ns	
Enzyme regulator	3.56%	2.74%	6.78%	3.82%	3.82% ns	7.54% *	
Signal transducer	0.69%	1.83%	2.01%	1.80%	1.91% ns	1.40% ns	

Table 2. The percentage of genes in each paralogy class belonging to each high-level GO functional category"

"Only the seven largest functional classes are shown, genes of unknown function have been removed. The total percentages can exceed 100% because a gene can be annotated as belonging to more than one class.

*Significance levels are ** P<0.001, * P<0.05 and not significant (ns) compared with the set of singleton genes.</p>



11 ... : 0



Orthologs vs. paralogs





Muil



http://www.ncbi.nlm.nih.gov/Education/BLASTinfo/Orthology.html

Mail


Homologs are most commonly defined as orthologs, paralogs, or xenologs.

> Orthologs are homologs produced by speciation—they represent genes derived from a common ancestor that diverged because of divergence of the organism. Orthologs tend to have similar function.

Paralogs are homologs produced by gene duplication and represent genes derived from a common ancestral gene that duplicated within an organism and then diverged. Paralogs tend to have different functions.

Xenologs are homologs resulting from the horizontal transfer of a gene between two organisms. The function of xenologs can be variable, depending on how significant the change in context was for the horizontally moving gene. In general, though, the function tends to be similar.

http://www.ncbi.nlm.nih.gov/About/primer/phylo.html



L Université de Lausanne Département d'écologie et évolution



"Ortholog conjecture"



11 ... : 0

UNIL | Université de Lausanne Département d'écologie et évolution



This should not be an assumption It should be an hypothesis to test

ortholog conjecture: functional change after duplication







IL | Université de Lausanne Département d'écologie et évolution



Alternative hypothesis:

Uniform distribution of positive selection



Another way of putting it





Département d'écologie et évolution



Study design



Unil

UNIL | Université de Lausanne Département d'écologie et évolution





Three cartoon models of evolution

- DDC subfunctionalization (loss of function) after duplication only
- Neofunctionalization after duplication only
- Neofunctionalization with time, irrespective of duplication
- > Different predictions in different study designs?

UNIL'| Université de Lausanne Département d'écologie et évolution



Comparisons inside one

genome



Département d'écologie et évolution





II. Duplicate pairs vs. random pairs

<--> ∨

What can we see under different models of evolution?

I- All models: Differences between paralogs

II- *All models*: Paralogs more similar than random pairs, but not identical



ó



III- *All models*: Measure of retention bias, confused by evolution after duplication



L | Université de Lausanne Département d'écologie et évolution



Biased gene retention





Département d'écologie et évolution

Unil

Arabidopsis WGD genes:

Affymetrix data from 62 different conditions and tissues



Correlation coefficient (r)

expression divergence: r < 0.52:

57% of pairs of recent duplicates 73% of pairs of ancient duplicates r < 0:

15% of pairs of recent duplicates29% of pairs of ancient duplicates



JNIL | Úniversité de Lausanne Département d'écologie et évolution

11 ... 0







Département d'écologie et évolution



Two genomes: functional data



11...:0

UNIL | Université de Lausanne Département d'écologie et évolution



IV. Orthologs of duplicates vs. orthologs of singletons



VI. Divergence relative to outgroup



IV- All models. Measure of retention bias

VI-*DDC* subfunctionalization after duplication: Two paralogs different, complementary to full outgroup function

Neofunctionalisation after duplication or with time: One paralog similar to outgroup, one different

VII. Divergence relative to outgroup of duplicates vs. singletons



VII- DDC subfunctionalization after duplication: Two paralogs different, complementary to outgroup; singleton similar to outgroup *Neofunctionalisation after duplication*: One paralog similar to outgroup, one different; singleton similar to

outgroup

Neofunctionalisation with time: No specific prediction

et évolution



a 			level of expression: high
d divergence in at least 2 tissues	subfunctionalization 19 (61)	asymmetric 49 (109)	Fisher test p-value 0.01 (0.01)
e divergence in exactly 2 tissues	15 (44)	41 (93)	0.02 (0.003)



UNIL | Université de Lausanne Département d'écologie et évolution



© 2009 SI



Genes highly similar in expression are underrepresented in ortholog sets with recent human- or mouse-specific gene duplications. Histograms of expression correlation coefficients (R) across 16 tissues are shown for (A) 1325 one-to-one orthologs between human and mouse; (B) 163 ortholog sets with one mouse sequence and more than one human co-ortholog; (C) 135 ortholog sets with one human sequence and more than one mouse co-ortholog

© 2009 SIB



Two genomes: sequences



Mail

UNIL | Université de Lausanne Département d'écologie et évolution



V. Divergence between orthologs of duplicates vs. between orthologs of singletons



V- All models: Measure of retention bias

VI. Divergence relative to outgroup



VI-*DDC subfunctionalization after duplication*: No prediction relative to symmetry, relaxed purifying selection *Neofunctionalisation after duplication* or *with time*: Asymmetry between paralogs, positive selection (low power)

VII. Divergence relative to outgroup of duplicates vs. singletons



VII- Higher divergence of duplicates, confused by retention bias



UL' Université de Lausanne Département d'écologie et évolution



After WGD in yeast: 76 / 457 gene pairs (17%) show accelerated protein evolution *defined as instances in which the amino acid substitution rate along one or both of the S. cerevisiae branches was at least 50% faster than the rate along the K. waltii branch.*







© 2009 SIB Zdobnov & Bork 2007 Trends Genet 23, 16-20

Département d'écologie et évolution

Real phylogenomics



Mmil

UNIL | Université de Lausanne Département d'écologie et évolution



VIII. Comparison of several singletons and duplicates per gene tree (sequences)



VII- DDC subfunctionalization after duplication: Higher relaxation of purifying selection on branches after duplication Neofunctionalisation after duplication: More positive selection on branches after duplication Neofunctionalisation with time: Positive selection in various branches of the tree

IX. Comparison of several singletons and duplicates per gene tree (functional data)



IX- DDC subfunctionalization after duplication: Conservation of pattern among singletons; sub-patterns in duplicates Neofunctionalisation after duplication:
Conservation in most homologs; new patterns in some duplicates Neofunctionalisation with time: Variation in pattern among homologs, with gain of new patterns





Discuss with your neighbor

- What is more frequent: point mutations or duplications?
- Where in animal evolution were there genome duplications?
- Two reasons why we should distinguish orthologs and paralogs
- What is gene function?
- Find an example of orthologs with different functions



Département d'écologie et évolution



Real data can be confusing



11 ... : 0

Département d'écologie





M/LWS pigments

shifts in wavelength

Zebrafish L-1(P558) Zebrafish L-2 (PS48) Cavefish L (PS58) Medaka L-A (P561) Medaka L-8 (P562) Cavefish M (P530) Clawed frog L (P557) Chameleon L (P560)* Human M (P530) Sq. monkey L (PSS8) Sq. monkey M (P545) Sq. monkey M (P532) Mouse M (P508)* Squirrel M (PS32) Dolphin M (P524)* Wallaby M (P528)

Yokoyama (2008) Ann Rev Genomics Human Genet 9: 259-282 Imil

UNIL | Université de Lausanne

et évolution

Département d'écologie





Tocchini-Valentini et al. 2009 J. Biol. Chem. 284: 1938-1948

SIB

Swiss Institute of Bioinformatics



Humans are not mice



Mmil

Département d'écologie et évolution



Review of human-mouse ortholog comparisons:

- ≥16% divergence of expression
- ≥13% divergence of splice forms
- 9% with secondary duplications
- ≥20% mutant lethal in human but not in mouse





11...:0

et évolution

Gharib and Robinson-Rechavi 2011 Briefings Bioinf 12, 436-441



Duplicates evolve asymmetrically



UNIL Université de Lausanne Département d'écologie et évolution





-> acceleration of one paralog



SIB

Swiss Institute of Bioinformatics

Brunet et al. (2006) Mol Biol Evol 23: 1808-16

et évolution





© 2009 SIB Conant & Wolfe 2008 Nature Rev Genet 9:938-950

UNIL Université de Lausanne Département d'écologie et évolution



Real neofunctionalization has postive selection



Département d'écologie et évolution



Species 1	Α	L	Р	н	Υ
	GCC	C <mark>T</mark> T	ССТ	CAT	TAT
Species 2	Α	R	Р	н	Υ
	GCC	C <u>G</u> T	ССТ	САТ	TA <u>C</u>

Measure of mutation, genetic drift and time:

• dS = <u>number of **synonymous** substitutions</u> synonymous sites

Measure of mutation, genetic drift and time and selection:

• dN = <u>number of **non-synonymous** substitutions</u> non-synonymous sites

Measure of selective pressure: ratio dN/dS (ω):

- dN/dS < 1 -> purifying selection
- $dN/dS = 1 \rightarrow neutral evolution$
- dN/dS > 1 -> positive selection

JNIL | Université de bausanne Département d'écologie et évolution



Codon models (PAML Package, Yang 2007)

- One ratio model
 - Same dN/dS for all branches and sites.
- Branch models
 - Estimate different dN/dS among branches.
- Site models
 - Estimate different dN/dS among sites.
- Branch-site models
 - Estimate different dN/dS among branches <u>and</u> sites.







NIL Université de bausanne Département d'écologie et évolution



Branch-site test



 $\omega = dN/dS$

 $\begin{array}{l} \omega_0 < 1 \\ \omega_1 = 1 \\ \text{proportions of sites: } p_0 + p_1 = 100\% \end{array}$

$$\omega_{0} < 1$$

$$\omega_{1} = 1$$

$$\omega_{2}$$

proportions of sites: $p'_{0} + p'_{1} + p'_{2} = 100\%$

$$\omega_{2} = 1$$

$$\omega_{2} \ge 1$$

background branches

Likelihood Ratio Test: Is the difference between the foreground and background branches due to positive selection?

© 2009 SIB

Département d'écologie et évolution

11 ... : 0




Bayesian Empirical Bayes posterior probability of positive selection for each site on this branch





Positive selection in vertebrate genes

Unil



Positive selection in vertebrate evolution



2673 branches tested (correction for multiple testing: q = 10%)

Positive selection significant for 45% of branches

Only 0.9% to 4.7% of sites

 \geq 1 branch significant for 77% of genes



Positive selection and duplication





Positive selection in vertebrate evolution



	Parameter	FSGD vs. other branches ^b		Euteleosts branch: FSGD topology vs. singleton		Bony vertebrates branch: 2R detected vs. not detected	
Branch behavior*		P-value ^c	Differenced	P-value ^e	Difference ⁴	P-value*	Difference ^d
LRT significant	ΔtnL	0.14	8.9-10.7	0.57	10-11	0.94	11-11
	Branch length*	0.017	0.085-0.11	0.014	0.093-0.12	0.035	0.13-0.14
	Mean wa	0.93	0.075-0.075	0.015	0.061-0.075	0.0098	0.072-0.079
	Percent of sites and	0.92	8396-8496	0.0029	89%-85%	0.67	83%-82%
	Percent of sites wy'	0.27	8.6%-10%	3.4×10^{-4}	6.6%-11%	0.035	9.3%-11%
	Percent of sites as,"	0.12	8.1%-5.6%	0.56	4.1%-4.1%	0.075	8,1%-7.2%
LRT nonsignificant	ΔLnL	0.41	0.69-0.65	0.30	0.79-0.65	0.33	0.91-1.0
	Branch length*	0.43	0.082-0.069	0.46	0.10-0.077	0.018	0.095-0.12
	Mean wo	0.13	0.054-0.062	0.65	0.061-0.061	3.8 × 10-6	0.048-0.064
	Percent of sites we	0.00026	80%-86%	0.28	84%-87%	0.011	88%-83%

Table 5. Influence of duplication on substitution parameters

Boldface indicates that the difference is significant.

*Classification according to whether the LRT for positive selection is significant on each branch (q = 0.10).

"Only branches from tree topologies for which the FSGD branch exists were used.

"Nonpaired Wilcoxon test. In bold if the difference is significant after Bonferroni correction (α = 0.05/30 = 0.0017).

^dDuplication or post-duplication mean - nonduplication mean.

"In amino acid substitutions/site.

Values for the foreground branch of each test.

(LRT) Likelihood ratio test.

• No impact on positive selection

• Relaxation of purifying selection

• Preferential retention of genes under strong purifying selection (Davis & Petrov 2004; Brunet et al 2006)

© 2009 SIB

11 ... : 0 Département d'écologie

et évolution



Amino acid shifts



Duplication and amino acid patterns



Gene expression



Unil





Database:



Bastian F., Parmentier G., Roux J., Moretti S., Laudet V., Robinson-Rechavi M., 2008. Lecture Notes in Computer Science 5109: 124-131

Mmil







Examples of gene expression query : zebrafish heart, adult zebrafish heart, adult zebrafish heart and adult human heart

Type 3 letters :			
Choose species ☑ Danio rerio ☑ Drosophila melanogaster ☑ Homo sapiens ☑ Mus musculus ☑ Xenopus tropicalis <u>Select All</u> <u>Ueselect All</u>	Choose stages Set All metastages Set embryo Set embryo Set embryo Set envage Set blastala Set gastrula Set organogenesis Set post-embryonic development Set adult Set embryonic development Set adult		
Choose anatomical structures alimentary system cardiovascular sy nervous system renal system reproductive system respiratory system skeletal system skeletal system	stem Selected anatomical structures : Search for ∉ heart ×		
Found 3 results for -kidn- kidn Name Part of Syn.: aduk kidney intraembryonic coelom, renal system, in Syn.: embryonic kidney intraembryonic coelom, renal system Syn.: metanephric kidney irenal system	amune vystem		

release 10: 15k Affymetrix chips 3k EST libraries

231k in situ experiments

⁺∕ SÎB

Swiss Institute of Bioinformatics

UNIL | Université de Lausanne Département d'écologie et évolution

Unil



Expression after duplication





Bastian F., Roux J., Robinson-Rechavi M., unpublished

11 ... : 0

UNIL | Université de Lausanne Département d'écologie et évolution



© 2009 SIB

Evolution of expression after duplication

• 4 types of "triplets":



An example: Khdrbs1 related genes

L	 -

Neofunctionalization

HOG name	Category	Expression singleton	Expression dupl. 1	Expression dupl. 2
cerebellum	Overlap	Expressed	Expressed	Expressed
archinephric duct	Overlap	Expressed	Expressed	Expressed
mesenchyme	Overlap	Expressed	Expressed	Expressed
somite	Overlap	Expressed	Expressed	Expressed
retina	Overlap	Expressed	Expressed	Expressed
ovary	Overlap	Expressed	Expressed	Expressed
testis	Overlap	Expressed	Expressed	Expressed
liver	Overlap	Expressed	Expressed	Expressed
beart	Overlap	Expressed	Expressed	Expressed
telencephalon	Overlap	Expressed	Expressed	Expressed
inner car	Overlap	Expressed	Expressed	Expressed
olfactory organ	Overlap	Expressed	Expressed	Expressed
viscerocranium	Overlap	Expressed	Expressed	Expressed
diencephalon	Overlap	Expressed	Expressed	Expressed
presomitic mesoderm	Overlap	Expressed	Expressed	Expressed
skin	Overlap	Expressed	Expressed	Expressed
tectum	Overlap	Expressed	Expressed	Expressed
forelimb - pectoral fin	Overlap	Expressed	Expressed	Expressed
mesonephros	Overlap	Not.expressed	Expressed	Expressed
pharynx	Overlap	Not expressed	Expressed	Expressed
seural tube	Neofunctionalization	Not extension	Expressed	Not expressed



UNIL | Université de Lausanne Département d'écologie et évolution



© 2009 SIB





SIB

Swiss Institute of Bioinformatics

Expression diverges (a bit) between paralogs



Département d'écologie et évolution



Let's work!

Unil



examine human gene you studied in Genomicus in:

http://selectome.unil.ch/

http://bgee.unil.ch/

- exclude Xenopus
- chose Data quality "high"
- ➢ show gene details
- display all homologous organ groups
- > overall picture of evolution of each of these genes



Département d'écologie et évolution



Support for alternative model



11 ... : 0



standard model: functional change after duplication



Département d'écologie et évolution



Time

Phylogenomics of function







© 2009 SIB



UniProtiCB Entry

PIR View

Nitegral View | SAS View

UniProtect Entry: P03372

ENTRY INFORMATION				
ENTRY NAME	ESRI-HUMAN			
ACCESSION NUMBERS	P03372; Q13511; Q14276; Q9NU51; Q9UDZ7; Q9UIS7			
CREATED	Release 01, 21-JUL-1986			
SEQUENCE UPDATE	Release 29,01-JUN-1994			
ANNOTATION UPDATE	Release 49, 24-JAN-2006			
NAME AND ORIGIN OF T	HE PROTEIN			
PROTEIN NAME Estrogen receptor				
DESCRIPTION	ER: Estradiol receptor: ER-alpha			
GENE NAME	ESR1; ESR; NR3A1			
SOURCE ORGANISM	Homo sapiens			
and the second se				

	CO-0016595 Casheematin namadaling complex NAS
	COMO10383 C. Chromaun remodeling complex JNAS.
	GO:0016020,C:membrane,NAS.
	GO:0030284 F:estrogen receptor activity,NAS.
	GO:0030235 F:nitric-oxide synthase regulator activity,NAS.
1	GO:0005515 F:protein binding JPI.
GO	GO:0016049 P:cell growth JSS.
1000	GO:0030520 P:estrogen receptor signaling pathway NAS.
	GO:0045839 P:negative regulation of mitosis ISS.
	GO:0006355 P:regulation of transcription, DNA-dependent, NAS.
	GO:0007165 P:signal transduction, TAS.
	OuickGO
-	

Unil UNIL | Université de Lausanne



Gene Ontology

- 17 genomes with highest GO coverage
- experimental evidence codes: EXP, IDA, IEP, IGI, IMP, IPI
- from GOA and Ensembl
- orthologs and paralogs inferred by OMA
- confirmation with Ensembl

Per species pair:

- one-to-one orthologs
- 1-to-many and many-to-many orthologs
- different-species paralogs.

Per species:

- same-species paralogs
- -> 26 330 pairs of genes with experimental GO annotations for both





Département d'écologie et évolution



Beware of the bias



Unil



Dataset	Same- Species Paralogs	Different- Species Paralogs	1:1 Orthologs	Other Orthologs
Same Publications	/ 1573	18	/ 154	44
Different Publications, Same Authors	613	382	874	434
No Common Author	2492	20719	13149	11296
All Experimental Annotations	3312	20766	13309	11371
/				
479	%	1%		



A. Authorship bias: average GO Similarity



C. Variation of background GO similarity among types of relations (random gene pairs)



B. Variation of GO term frequency among species







Comparisons after correcting for biases



Département d'écologie et évolution





Functional Similarity of Orthologs and Paralogs Among All 13 Species



et évolution

Orthologs are useful



Unil





Functional Similarity of corresponding 1:1 Orthologs in Human, Mouse and Outgroup

Unil



Conclusion orthologs/ paralogs



11 ... : 0



The standard model holds, but is weaker than expected.

Duplication does affect gene evolution Other factors are significant in gene function divergence





General conclusion

Mail


- Duplication happens:
 - Duplication is frequent
 - Genome duplication has shaped many eukaryotic genomes
- Function:
 - Duplicate genes can diverge and add new functions
 - Changes in function can concern
 - expression
 - protein sequence
 - interactions
 - and more
 - Duplicate genes can be conserved
 - Orthologs can diverge

Duplication is a major force in the evolution of genomes



NIL | Úniversité de Lausanne Département d'écologie et évolution





Unil

UNIL | Université de Lausanne Département d'écologie et évolution

