# Gradual divergence and diversification of mammalian duplicate gene functions

Raquel Assis[1]*, Doris Bachtrog[2]

[1] Department of Biology, Huck Institutes of Life Sciences, Pennsylvania State University, University Park, PA 16802.

[2] Department of Integrative Biology, Center for Theoretical Evolutionary Genomics, University of California, Berkeley, CA 94720, USA.

*Correspondence to: rassis@psu.edu

Running Title: Functional divergence of mammalian duplicate genes

Keywords: Gene duplication, neofunctionalization, subfunctionalization, gene expression divergence

1    **Gene duplication provides raw material for the evolution of functional innovation. We recently**

2    **developed a phylogenetic method to classify the evolutionary processes underlying the retention and**

3    **functional evolution of duplicate genes by quantifying divergence of their gene expression profiles.**

4    **Here, we apply our method to pairs of duplicate genes in eight mammalian genomes, using data from**

5    **11 distinct tissues to construct spatial gene expression profiles. We find that young mammalian**

6    **duplicates are often functionally conserved, and that functional divergence gradually increases with**

7    **evolutionary distance between species. Examination of expression patterns in genes with conserved**

8    **and new functions supports the "out-of-testes" hypothesis, in which new genes arise with testis-**

9    **specific functions and acquire functions in other tissues over time. While new functions tend to be**

10   **tissue-specific, there is no bias toward expression in any particular tissue. Thus, duplicate genes**

11   **acquire a diversity of functions outside of the testes, possibly contributing to the origin of a multitude**

12   **of complex phenotypes during mammalian evolution.**

13

14   **Introduction**

15   Gene duplication produces copies of existing genes, which can diverge from their ancestral states and

16   contribute to the evolution of novel phenotypes. A large proportion of mammalian genes arose via gene

17   duplication (Li et al. 2001; Ryvkin et al. 2009), and many are members of gene families with diverse and

18   essential functions. For example, Hox, growth factor, and olfactory receptor gene families were all

19   produced by gene duplication.  However, the evolutionary paths leading from functionally redundant

20   duplicate copies to distinct genes with important functions remain unclear.

21

22   Different processes may drive the long-term retention and functional evolution of duplicate genes:

23   Parent and child copies may each maintain the function of their ancestral single-copy ortholog

24   (conservation; Ohno 1970); one copy may maintain the ancestral function, while the other acquires a

1    new function (neofunctionalization; Ohno 1970); each copy may lose part of its function, such that

2    together both copies carry out the ancestral function (subfunctionalization; Force et al. 1999; Stoltzfus

3    1999); or both copies may acquire new functions (specialization; He and Zhang 2005). We recently

4    developed a method that utilizes distances between gene expression profiles to classify these

5    evolutionary processes (Assis and Bachtrog 2013). Our method is applied to pairs of duplicates and

6    requires that, for each pair, we can distinguish between parent and child copies and identify a single-

7    copy ancestral ortholog in a closely related sister species. Moreover, parent, child, and ancestral genes

8    must all have spatial or temporal gene expression data from which gene expression profiles can be

9    constructed.

10

11    To study the roles of conservation, neofunctionalization, subfunctionalization, and specialization in the

12    retention of mammalian duplicate genes, we applied our method to pairs of duplicate genes from eight

13    mammalian genomes: human (*Homo sapiens*), chimpanzee (*Pan trogodytes*), gorilla (*Gorilla gorilla*),

14    orangutan (*Pongo pygmaeus abelii*), macaque (*Macaca mulatta*), mouse (*Mus musculus*), opossum

15    (*Monodelphis domestica*), and platypus (*Ornithorhynchus anatinus*). Using synteny information from

16    whole-genome alignments to determine orthologous genomic positions, and parsimony to infer the

17    evolutionary dynamics of genes, we distinguished between parent and child copies and identified

18    ancestral single-copy orthologs for each pair of duplicates. Then, we applied our classification method to

19    RNA-seq data from 11 mammalian tissues: female and male cerebrum, female and male cerebellum,

20    female and male heart, female and male kidney, female and male liver, and testis (Brawand et al. 2011).

21

22    **Results**

23    In total, we obtained 654 pairs of mammalian duplicate genes for which we could distinguish between

24    parent and child copies and also identify at least one expressed single-copy ancestral gene in a closely

3

1    related sister species. Application of our method to these pairs yielded 382 cases of conservation, 213

2    cases of neofunctionalization (105 neofunctionalized parent copies, 108 neofunctionalized child copies),

3    9 cases of subfunctionalization, and 50 cases of specialization. Thus, most mammalian duplicate genes

4    have conserved functions. Moreover, functional divergence typically affects only one gene copy, and

5    retention of duplicates by subfunctionalization is rare.

6

7    Comparing duplicates from mammalian genomes of different evolutionary distances enabled us to

8    examine whether there is a negative relationship between functional conservation and age of duplicate

9    genes, as expected if genes evolve new functions over time. We used parsimony to date the origin of

10   child copies along the mammalian phylogeny (Figure 1A). Consistent with global patterns, conservation

11   is the most common evolutionary process underlying the retention of duplicate genes in every

12   mammalian lineage examined (Figure 1A). To test if functional conservation decreases with increasing

13   evolutionary divergence between species, we calculated rates of protein sequence divergence ($K_a$)

14   between single-copy genes in human and each sister species on the tree, and used these values as

15   estimates of evolutionary divergence between pairs of species. Comparison of median $K_a$ to proportions

16   of duplicate gene pairs with conserved functions revealed that functional conservation of duplicates

17   indeed decreases significantly with evolutionary divergence between species, and that this decrease is

18   approximately linear (Figure 1B). Thus, young pairs of mammalian duplicates are generally functionally

19   conserved, and new functions evolve gradually over time. Moreover, $K$a is not as strongly correlated to

20   the proportion of functionally conserved single-copy genes, indicating that functional divergence occurs

21   faster in duplicates than in single-copy genes.

22

23   To determine the types of novel functions acquired by mammalian duplicates over time, we examined

24   differences between gene expression patterns in copies of pairs retained by neofunctionalization. In

1    such cases, one copy has maintained the ancestral function (the "conserved" copy), while the other has

2    acquired a new function (the "neofunctionalized" copy). Thus, we can directly assess ancestral and new

3    functions within pairs. We used the highest relative tissue expression level for each gene as a measure

4    of its tissue specificity. Comparison of distributions of tissue specificities revealed that, as expected,

5    conserved copies and ancestral genes have similar tissue-specific expression levels (Figure 2A). Single-

6    copy genes have similar tissue-specific levels as well, indicating that duplicate genes are initially as

7    broadly expressed as single-copy genes in the genome. In contrast, neofunctionalized copies are

8    significantly more tissue-specific than ancestral, conserved, and single-copy genes. An alternative metric

9    of tissue-specific expression, τ (Yanai et al. 2005), yields the same conclusions (Figure S1). Thus,

10   duplicates are initially as broadly expressed as typical single-copy genes, and acquisition of a new

11   function by neofunctionalization results in increased tissue specificity.

12

13   To assess whether conserved and neofunctionalized copies of pairs show different expression patterns

14   across tissues, we examined quantities of genes with highest expression levels in each tissue (Figure 2B).

15   In most tissues, numbers of highly expressed single-copy genes are similar to those of conserved gene

16   copies. The two exceptions are male liver and testis. While the difference in highly expressed male liver

17   genes is modest, the proportion of conserved testis-specific genes is nearly double that of single-copy

18   testis-specific genes. In contrast, only a small proportion of neofunctionalized copies are testis-specific.

19   Thus, many genes initially arise with testis-specific functions and gradually acquire other functions over

20   time, supporting the "out-of-testes" hypothesis of new gene origination (Kaessmann 2010). Moreover,

21   while functions of neofunctionalized copies are typically more tissue-specific (Figure 2A), there is no bias

22   toward specificity in any particular tissue(s). Hence, it appears that mammalian duplicate genes acquire

23   new functions in a diversity of tissues.

24

5

1    **Discussion**

2    In a recent study, we applied our classification method to pairs of duplicate genes in *Drosophila*

3    *melanogaster* and *D. pseudoobscura* (Assis and Bachtrog 2013), for which the median $K$s is 1.79

4    (Richards et al. 2005). Contrary to our observation in mammalian duplicates, we found that most

5    *Drosophila* duplicates were neofunctionalized, and examination of evolutionary processes over shorter

6    divergence times suggested that novel functions arise rapidly (Assis and Bachtrog 2013). However, the

7    smallest $K$s examined in *Drosophila* was 0.11 (between *D. melanogaster* and *D. simulans*) (Lazzaro 2005).

8    In contrast, while the $K$s between human and platypus is approximately 1.41 (Warren et al. 2008), the

9    smallest $K$s examined in mammals was 0.01 (between human and chimpanzee; Chen and Li 2001), which

10   is an order of magnitude smaller than that between any pair of *Drosophila* species we analyzed (Assis

11   and Bachtrog 2013). Thus, we have greater temporal resolution in mammals than in *Drosophila*,

12   enabling us to more closely examine the functional diversification of mammalian duplicates over

13   evolutionary time.

14

15   In contrast to the widespread and rapid neofunctionalization observed in *Drosophila*, young mammalian

16   duplicates are primarily conserved, and new functions arise slowly over time. This difference may be due

17   to the larger effective population size ($N_e$) of *Drosophila* than of mammals (Beckenbach et al. 1993;

18   Lynch and Conery 2003; Jensen and Bachtrog 2011), which contributes to more efficient adaptive

19   protein sequence evolution in *Drosophila* (Britten 1986; Moriyama 1987; Carroll 2005), and could

20   similarly result in more rapid acquisition of adaptive functions by *Drosophila* duplicate genes.

21   Furthermore, while small $N_e$ is also thought to result in a higher prevalence of subfunctionalization

22   (Lynch et al. 2001), this process does not appear to play a major role in the retention of duplicate genes

23   in either lineage. A possible reason for this observation is that subfunctionalization may be more

6

1  common in duplicate genes produced by whole genome duplication events (Casneuf et al. 2006; Fares et

2  al. 2013), which our study does not examine.

3

4  In both *Drosophila* and mammalian duplicate genes, we uncovered strong support for the "out-of-

5  testes" hypothesis of new gene emergence (Kaessmann 2010). Testes may facilitate the initial

6  transcription of young genes, while sheltering them from pseudogenization as they acquire new

7  functions (Kaessmann 2010), and may thus be an ideal tissue for young genes. However, while there

8  does not appear to be a bias in evolution of functions from testes to any particular tissue in either

9  lineage, a major difference we observe is that functions of *Drosophila* genes broaden over time, whereas

10  functions of mammalian genes narrow over time. Hence, *Drosophila* duplicates eventually acquire broad

11  housekeeping functions, whereas mammalian duplicates acquire tissue-specific functions that may

12  facilitate the evolution of phenotypic diversity across species.

13

14  **Methods**

15  ***Identification of duplicate and single-copy genes***

16  We downloaded protein sequences, annotation files, and lists of duplicate genes for all genomes from

17  the Ensembl database at http://www.ensembl.org. To obtain a comprehensive list of duplicates in each

18  mammalian genome, we supplemented Ensembl lists with those from the Duplicated Genes Database

19  (DGD) at http://www.dgd.genouest.org and with protein BLAST searches (Altschul et al. 1990), which we

20  performed as previously described (Assis and Bachtrog 2013). Any annotated genes not on these lists

21  were considered to be single-copy genes, and gene families with more than two copies were excluded

22  from our analysis. We quantile-normalized RNA-seq data from mammalian tissues (Brawand et al. 2011)

23  and restricted our analysis to pairs for which both copies are expressed in at least one tissue. Using

1     these expression data, we classified pairs of duplicates as conserved, neofunctionalized,

2     subfunctionalized, or specialized as previously described (Assis and Bachtrog 2013).

3

4     ***Phylogenetic dating and identification of ancestral single-copy orthologs***

5     We downloaded whole-genome alignments from Ensembl (http://www.ensembl.org) and UCSC Genome

6     Bioinformatics (http://www.genome.ucsc.edu) databases and extracted syntenic regions in all genomes

7     for each duplicate gene. We used parsimony to phylogenetically date the origin of each pair of

8     duplicates. Duplicates that were present in all species or that could not be resolved via parsimony were

9     removed from our analysis. For each pair, the gene copy aligning to the ancestral single-copy gene(s)

10     was considered the parent, and the second copy was considered the child. Orthologs for single-copy

11     genes were also obtained via synteny and aligned with MACSE (Ranwez et al. 2011). PAML (Yang 2007)

12     was used to estimate the $K_a$ between each pair of single-copy genes.

13

14     ***Statistical analyses***

15     Mann-Whitney *U* tests were used to compare distributions of relative expression levels, and Fisher's

16     Exact tests were used to compare observed and expected numbers of genes with highest relative

17     expression in each tissue for conserved and neofunctionalized classes, as well as to compare numbers of

18     genes between classes. All statistical analyses were performed in the R software environment (R
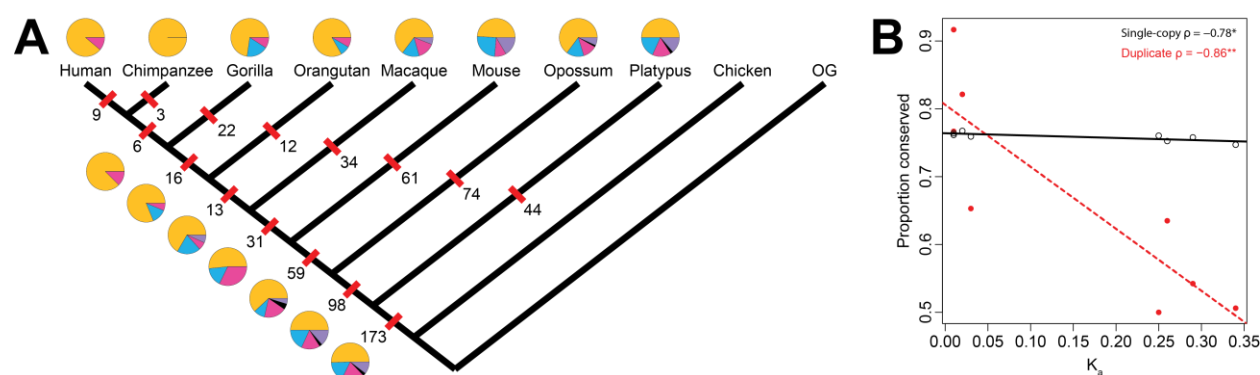
19     Development Core Team 2009).

20

21

22

23

24

**Figure 1**. **Evolutionary processes driving the retention of mammalian duplicate genes.** *A*) Pie charts depicting the role of each process on different branches of the mammalian phylogeny (yellow = conserved; blue = neofunctionalization of parent copy; pink = neofunctionalization of child copy; black = subfunctionalization; purple = specialization). Numbers of duplicate gene pairs examined along each branch are indicated beside red ticks. *B*) Relationship of median $K_a$ between pairs of species to proportions of functionally conserved single-copy genes (black) and pairs of duplicate genes (red). Linear regression lines are depicted to show rate of decreased functional conservation, and Pearson's correlation coefficients are shown in the top right corner of the plot. * $p < 0.05$; ** $p < 0.01$.
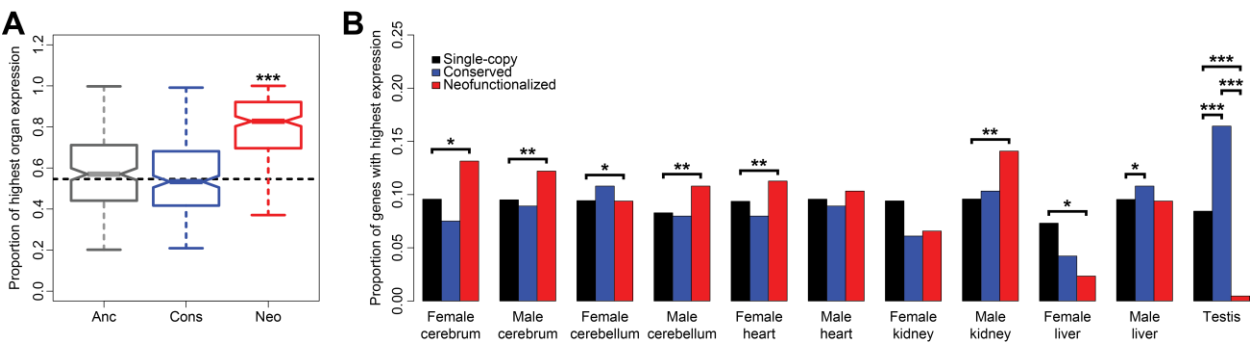
**Figure 2. Comparison of tissue-specific expression in conserved and neofunctionalized copies of pairs that underwent neofunctionalization.** *A*) Boxplots of highest relative expression levels for ancestral (Anc, gray), conserved (Cons, blue), and neofunctionalized (Neo, red) genes. Dotted black line represents median for single-copy genes, and asterisks show significance relative to distribution of single-copy genes. *B*) Barplots depicting proportions of single-copy (black), conserved (blue) and neofunctionalized (red) genes with highest expression in each tissue. Asterisks above lines connecting two bars indicate significance between classes. * $p < 0.05$; ** $p < 0.01$; $p < 0.001$.
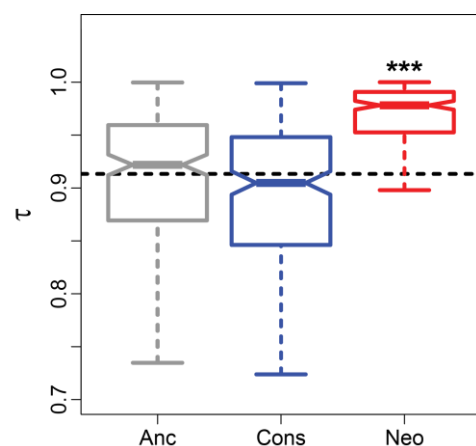
1

2 **Figure S1. Comparison of tissue-specific expression in conserved and neofunctionalized copies of pairs**

3 **that underwent neofunctionalization.** Boxplots of tissue specificity indices ($\tau$) for ancestral (Anc, gray),

4 conserved (Cons, blue), and neofunctionalized (Neo, red) genes. Dotted black line represents median for

5 single-copy genes, and asterisks show significance relative to distribution of single-copy genes.

6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29

**References**

Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J Mol Biol* 215: 403-410.

Assis R, Bachtrog D. 2013. Neofunctionalization of young duplicate genes in *Drosophila*. *Proc Natl Acad Sci USA* **110:** 17409-17414.

Beckenbach AT, Wei YW, Liu H. 1993. Relationships in the *Drosophila obscura* species group inferred from mitochondrial cytochrome oxidase II sequences. *Mol Biol Evol* **10:** 619-634.

Brawand D, Soumillon M, Necsulea A, Julien P, Csardi G, et al. 2011. The evolution of gene expression levels in mammalian organs. *Nature* **478:** 343-348.

Britten RJ. 1986. Rates of DNA sequence evolution differ between taxonomic groups. *Science* **231:** 1393-1398.

Carroll SB. 2005. Evolution at two levels, on genes and form. *PLoS Biol* **3:** e245.

Casneuf T, Bodt SD, Raes J, Maere S, Van de Peer Y. 2006. Nonrandom divergence of gene expression following gene and genome duplications in the flowering plant *Arabidopsis thaliana*. *Genome Biol* **7:** R13.

Chen FC, Li WH. 2001. Genomic divergences between humans and other hominoids and the effective population size of the common ancestor of humans and chimpanzees. *Am J Hum Genet* **68:** 444-456.

Fares MA, Keane OM, Toft C, Carretero-Paulet L, Jones GW. 2013. The roles of whole-genome and small-scale duplications in the functional specialization of *Saccharomyces cerevisiae* genes. *PLoS Genet* **9:** e1003176.

Force A, Lynch M, Pickett FB, Amores A, Yan Y, et al. 1999. Preservation of duplicate genes by complementary, degenerative mutations. *Genetics* **151:** 1531-1545.

He X, Zhang J. 2005. Rapid subfunctionalization accompanied by prolonged and substantial neofunctionalization in duplicate gene evolution. *Genetics* **169:** 1157-1164.

Jensen JD, Bachtrog D. 2011. Characterizing the influence of effective population size on the rate of adaptation, Gillespie's Darwin domain. *Genome Biol Evol* **3:** 687-701.

Kaessmann H. 2010. Origins, evolution, and phenotypic impact of new genes. *Genome Res* **20:** 1313-1326.

Lazzaro B. 2005. Elevated polymorphism and divergence in the class c scavenger receptors of *Drosophila melanogaster* and *D. simulans*. *Genetics* **169:** 2023-2034.

Li WH, Gu Z, Wang H, Nekrutenko A. 2001. Evolutionary analyses of the human genome. *Nature* **409:** 847-849.

12

1   Lynch M, Conery S. 2003. The origins of genome complexity. *Science* **302:** 1401-1404.
2
3   Lynch M, O'hely M, Walsch B, Force A. 2001. The probability of preservation of a newly arisen gene
4   duplicate. *Genetics* **159:** 1789-1804.
5
6   Moriyama EN. 1987. Higher rates of nucleotide substitution in *Drosophila* than in mammals. *Jpn J*
7   *Genetics* **62:** 139-147.
8
9   Ohno S. 1970. Evolution by gene duplication. Springer-Verlag, Berlin.
10
11  R Development Core Team. 2009. R, A Language and Environment for Statistical Computing. R
12  Foundation for Statistical Computing. Vienna, Austria.
13
14  Ranwez V, Harispe S, Delsuc F, Douzery EJP. 2011. MACSE, Multiple Alignment of Coding SEquences
15  accounting for frameshifts and stop codons. *PLoS ONE* **6:** e22594.
16
17  Richards S, Liu Y, Bettencourt BR, Hradecky P, Letovsky S, et al. 2005. Comparative genome sequencing
18  of *Drosophila pseudoobscura*: chromosomal, gene, and cis-element evolution. *Genome Res* **15:** 1-18.
19
20  Ryvkin P, Jun J, Hemphill E, Nelson C. 2009. Duplication mechanisms and disruptions in flanking regions
21  influence the fate of mammalian gene duplicates. *J Comput Biol* **16:** 1253-1266.
22
23  Stoltzfus A. 1999. On the possibility of constructive neutral evolution. *J Mol Evol* **49:** 169-181.
24
25  Warren WC, Hillier LW, Marshall Graves JA, Birney E, Ponting CP, et al. 2008. Genome analysis of the
26  platypus reveals unique signatures of evolution. *Nature* **453:** 175-183.
27
28  Yanai I, Benjamin H, Shmoish M, Chalifa-Caspi V, Shklar M, et al. 2005. Genome-wide midrange
29  transcription profiles reveal expression level relationships in human tissue specifications. *Bioinformatics*
30  **21:** 650–659.
31
32  Yang Z. 2007. PAML 4, Phylogenetic Analysis by Maximum Likelihood. *Mol Biol Evol* **24:** 1586-1591.
33
34
35
36
37
38
39
40
41