



## Unveiling the functional fate of duplicated genes through expression profiling and structural analysis

Alex Warwick Vesztrocy, Natasha Glover, Paul D. Thomas, et al.

*Genome Res.* 2025 35: 2273-2284 originally published online August 14, 2025

Access the most recent version at doi:[10.1101/gr.280166.124](https://doi.org/10.1101/gr.280166.124)

---

**References** This article cites 79 articles, 17 of which can be accessed free at:  
<http://genome.cshlp.org/content/35/10/2273.full.html#ref-list-1>

**Creative Commons License** This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <https://genome.cshlp.org/site/misc/terms.xhtml>). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

**Email Alerting Service** Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).

---

An advertisement banner with a teal background. On the left, the text reads 'CRISPR and RNAi Genetic Screening. Your new superpower.' In the center, there is a white-bordered box containing the words 'LEARN MORE'. On the right, there is a photograph of a woman wearing a red mask and a red cape, and the Cellecta logo, which consists of a green molecular structure and the word 'CELLECTA' in white capital letters.

CRISPR and RNAi Genetic Screening.  
Your new superpower.

LEARN MORE

CELLECTA

---

To subscribe to *Genome Research* go to:  
<https://genome.cshlp.org/subscriptions>

---

© 2025 Warwick Vesztrocy et al.; Published by Cold Spring Harbor Laboratory Press

## Method

# Unveiling the functional fate of duplicated genes through expression profiling and structural analysis

Alex Warwick Vesztröcy,<sup>1,2,3</sup> Natasha Glover,<sup>1,2</sup> Paul D. Thomas,<sup>4</sup> Christophe Dessimoz,<sup>1,2</sup> and Irene Julca<sup>1,2</sup>

<sup>1</sup>SIB Swiss Institute of Bioinformatics, 1015 Lausanne, Switzerland; <sup>2</sup>Department of Computational Biology, University of Lausanne, 1015 Lausanne, Switzerland; <sup>3</sup>BioSoft Research UK, London EC2A 4NE, United Kingdom; <sup>4</sup>Department of Population and Public Health Sciences, Keck School of Medicine, University of Southern California, Los Angeles, California 90033, USA

Gene duplication is a major evolutionary source of functional innovation. Following duplication events, gene copies (paralogs) may undergo various fates, including retention with functional modifications (such as subfunctionalization or neofunctionalization) or loss. When paralogs are retained, this results in complex orthology relationships, including one-to-many or many-to-many. In such cases, determining which one-to-one pair is more likely to have conserved functions can be challenging. It has been proposed that, following gene duplication, the copy that diverges more slowly in sequence is more likely to maintain the ancestral function—referred to here as “the least diverged ortholog (LDO) conjecture.” This study explores this conjecture, using a novel method to identify asymmetric evolution of paralogs and applying it to all gene families across the Tree of Life in the PANTHER database. Structural data for over 1 million proteins and expression data for 16 animals and 20 plants are used to investigate functional divergence following duplication. This analysis, the most comprehensive to date, reveals that, whereas the majority of paralogs display similar rates of sequence evolution, significant differences in branch lengths following gene duplication can be correlated with functional divergence. Overall, the results support the least diverged ortholog conjecture, suggesting that the least diverged ortholog tends to retain the ancestral function, whereas the most diverged ortholog (MDO) may acquire a new, potentially specialized role.

[Supplemental material is available for this article.]

Gene duplication is widely recognized as a key mechanism driving genome evolution (Ohno 1970). Indeed, redundant gene copies provide a “temporary escape from the relentless pressure of natural selection” (Ohno 1972), allowing the evolution of new functions. These genes that result from duplication events are referred to as “paralogs,” in contrast to “orthologs,” which descend from a common ancestral gene (Fitch 1970).

The role of duplication in the emergence of new gene functions has been extensively discussed (Hurles 2004). It is generally accepted that orthologs tend to be more conserved in function compared to paralogs—a concept known as the “ortholog conjecture” (Nehrt et al. 2011; Altenhoff et al. 2012; Gabaldón and Koonin 2013). Following a duplication event, the two gene copies can be retained with functional redundancy, diverge (sub- or neofunctionalization), or become lost (Force et al. 1999; Sémon and Wolfe 2007; Leitch and Leitch 2008). The most common phenomenon is that one copy accumulates deleterious mutations and thereby is silenced and eventually pseudogenized. Alternatively, mutations can lead to functional divergence between the duplicates, resulting in subfunctionalization (each copy may retain part of the ancestral function), or neofunctionalization (one copy may retain the ancestral function, whereas the other acquires a new function) (Lynch and Force 2000; Adams and Wendel 2005; Sehrish et al. 2014). Thus, gene duplication can provide the raw material for the evolution of functional innovation.

Asymmetric sequence evolution of paralogs can correlate with asymmetric functional divergence. It is generally accepted

that, following duplication, the two gene copies typically undergo a period of accelerated evolution (Huminićki and Wolfe 2004), which may be apparent only in one copy (Huerta-Cepas et al. 2011; Pegueroles et al. 2013; Pich I Roselló and Kondrashov 2014). This asymmetry has been interpreted as supporting the Ohno model of evolution (Scannell and Wolfe 2008; Pegueroles et al. 2013), which hypothesizes that one copy (“slow” copy, shorter branch in a gene tree) maintains the ancestral rate of evolution and function, whereas the other copy (“fast” copy, longer branch) may acquire a novel function (Ohno 1970).

The ortholog conjecture has motivated the notion of a “least diverged ortholog” (LDO) in the PANTHER database of phylogenetic trees (Mi et al. 2010). To define subfamilies, branch lengths are compared for each lineage following duplication, and only the subtree with the shortest branch is retained as the LDO subfamily. In this way, a single orthologous pair can be chosen even in the presence of postspeciation duplication events, with the pair being expected to be more likely than other pairs to retain the ancestral function. The advantage of this approach is that the least diverged orthologs are distinguishable in all cases, except those where the branch lengths following duplication are identical. The disadvantage, however, is that even small amounts of stochastic variation in the distance estimation will be conflated with actual differences in evolutionary rates.

Here, a generalization of this approach is used: in instances where gene duplications arise following a speciation event,

**Corresponding author: irene.julca@unil.ch**

Article published online before print. Article, supplemental material, and publication date are at <https://www.genome.org/cgi/doi/10.1101/gr.280166.124>.

© 2025 Warwick Vesztröcy et al. This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <https://genome.cshlp.org/site/misc/terms.xhtml>). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

orthologous relationships can take the form of one-to-many or many-to-many. In this study, the copy with the shorter branch in the gene tree is termed the “least diverged ortholog,” and the copy with the longer branch is called the “most diverged ortholog” (MDO). Furthermore, the notion in PANTHER that LDO gene pairs are more likely to retain the ancestral function can also be applied to this framework: after duplication, the gene copy with the shorter branch (LDO) is more likely to maintain the ancestral function than that with the longer branch (MDO). Although previous work did not use this terminology, the notion that the LDO tends to retain ancestral function remains an open question. This is referred to here as the “Least Diverged Ortholog Conjecture.”

Several studies have investigated aspects of the LDO conjecture empirically. Some mainly focused on the “ortholog conjecture,” whereas others investigated the asymmetrical evolution of paralogs. These studies have identified a distinct rate of sequence evolution between paralogs, indicative of functional differences (Huminięcki and Wolfe 2004; Cusack and Wolfe 2007; Scannell and Wolfe 2008; Pegueroles et al. 2013; Pich I Roselló and Kondrashov 2014). To assess functional divergence between paralogous genes, various data sets have been employed, including Gene Ontology (GO) terms (Nehrt et al. 2011; Altenhoff et al. 2012), protein-protein interactions (Bandyopadhyay et al. 2006; Kim and Yi 2006), and expression patterns (Gu et al. 2005; Li et al. 2005; Wang et al. 2012; Assis and Bachtrog 2013). Additionally, the fast-evolving copy has been associated with increased levels of tissue specificity (Scannell and Wolfe 2008; Huerta-Cepas et al. 2011). Furthermore, some studies have demonstrated that even the slower evolving copy shows evidence of a burst of protein sequence evolution immediately after duplication (Scannell and Wolfe 2008), whereas others have observed an overall increase in evolutionary rate postduplication not typically associated with rate asymmetry (Vance et al. 2022). However, the majority of these studies have relied on sequence differences ( $K_a/K_s$ ,  $d_N/d_S$ ) to differentiate between the fast- and slow-evolving copies, along with expression data from a restricted set of species and anatomical samples.

To investigate the LDO conjecture, this study aims to develop a novel method for capturing differences in selective pressures by identifying significant rate shifts in branch lengths following gene duplication events. Additionally, it seeks to explore functional differences between evolutionary models with statistically significant rate differences by analyzing protein structural data for over 1 million proteins and integrating mRNA expression data from 16 animal and 20 plant species. Finally, the study examines the specialization hypothesis, which proposes that the least diverged ortholog typically retains ancestral function, whereas the other copy evolves into a more divergent and specialized role.

## Results

### Least diverged ortholog conjecture analysis across the Tree of Life

This study investigated the least diverged ortholog conjecture, which implies that the gene copy with the shortest branch in the gene tree tends to retain the ancestral function after a duplication event. Analysis was performed on a data set comprising 143 organisms representing different lineages of the Tree of Life, across 15,693 gene trees. This conjecture relies solely on observed branch length and classifies gene copies with the shorter branch as the least diverged ortholog and those with the longer branch as the most diverged ortholog.

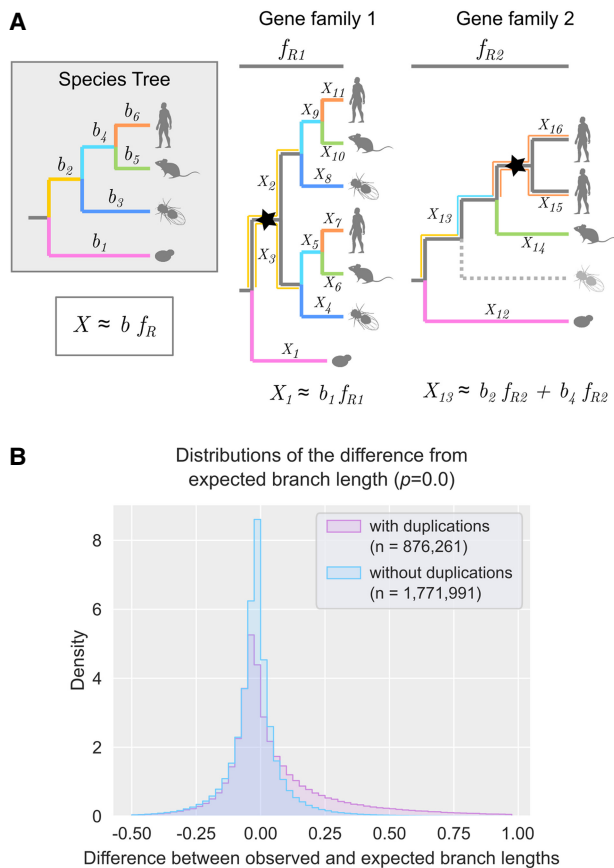
However, when gene families have different evolutionary rates, it is not clear whether small differences in branch length estimates can accurately classify gene copies into these two categories (LDO/MDO). Here, in order to detect distinct evolutionary rates on each of the branches, a statistical method was developed that allows for these family-specific evolutionary rates. This method is able to identify branches in a gene tree with accelerated rates of protein sequence evolution using only branch length information. This offers two key advantages: (1) it allows us to leverage pre-compiled libraries of gene trees such as PANTHER, eliminating the need to construct a custom data set; and (2) it enables the analysis of ancient gene duplication events, where nucleotide substitutions have reached saturation and the usual approach of using  $K_a/K_s$  cannot be reliably applied (Smith and Smith 1996; Gharib and Robinson-Rechavi 2013). Additionally, this method assigns an expected branch length to each gene tree branch, providing a model-based score that reflects deviations from the typical evolutionary rate of the gene family. As a result, this method attempts to classify a branch as typical, short, or long relative to the expected rate for the gene family.

To capture differences in evolutionary rates across branches in each gene tree, the method estimates gene family-specific rates of evolution (family-specific factors) for each gene tree (Fig. 1A). A factor greater than one indicates that a gene family is evolving more rapidly than average, potentially reflecting relaxed functional constraints, whereas a factor below one suggests slower evolution, which may reflect stronger functional constraints. The distribution of these factors has a mean of 1.06 and a standard deviation of 0.76 (Supplemental Fig. S1; Supplemental Table S1). A slight skew toward gene families with factors below one suggests that a considerable number may be evolving under stronger evolutionary constraints.

Using the family-specific factors, the expected branch length was computed for all branches in every gene tree (Fig. 1A) before computing the Z-score of the differences between the observed and expected branch lengths. A significant increase in evolutionary rate was observed between branches with no duplication (average  $-0.02$ ) and those with duplications ( $0.07$ ;  $P < 0.001$  Mann–Whitney  $U$  test) (Fig. 1B), suggesting that paralogs tend to evolve faster than orthologs. Similar results were observed when each major lineage was analyzed independently (Supplemental Fig. S2).

Following this analysis, duplication events were identified across all gene trees and classified into six categories (Fig. 2A; Supplemental Table S2): (i) normal-normal: both branches have a typical rate of evolution and are consistent with their expected lengths; (ii) short-short: both branches are significantly shorter than expected; (iii) long-long: both branches are significantly longer than expected; (iv) normal-long: only one branch is significantly longer than expected; (v) short-normal: only one branch is significantly shorter than expected; and (vi) short-long: one branch is significantly shorter than expected and the other is significantly longer than expected.

The analysis gave a total of 926,814 duplication events across the Tree of Life, with an average of 68 duplications per family. When examining the major lineages (Fig. 2B, gray bars), the Viridiplantae lineage exhibits the highest percentage of duplications from the total number of duplications (68%) (Fig. 2B; Supplemental Table S3), followed by Deuterostomia (16%) and Protostomia (4%). Internal nodes within this tree contributed only 5% of the total number of duplications (Supplemental Table S3). These findings highlight the significant role of gene



**Figure 1.** Analysis and calculation of expected gene tree branch lengths. (A) Graphical illustration of the calculation of the expected branch lengths ( $X$ ). On the left, a species tree depicts the phylogenetic relationships of four species with branches highlighted by different colors and branch numbers ( $b_1$  to  $b_6$ ). On the right, two gene trees show a duplication event (node marked with a star). The first gene tree is complete, and the second tree has a gene loss represented by a dashed line.  $f_R$  denotes the family-specific rate. When no losses occur (Gene family 1), there is an exact correspondence between the speciation branch in the gene tree and the species tree ( $X_1$  to  $X_{11}$ ). When there is a loss (Gene family 2), the speciation branch in the gene tree represents the sum of the corresponding branches in the species tree ( $X_{13}$ ). Note that in both gene families, when a duplication occurs, all paralogous branches are calculated, and the expected values remain the same. (B) Distributions of the difference from expected branch lengths in branches with duplications (pink) and without duplications (light blue). The sample size ( $n$ ) is indicated in the figure. The  $P$ -value for the Mann–Whitney  $U$  test is shown between parentheses ( $P=0.0$ ).

duplication in the evolutionary history of plants, in agreement with previous studies (Jiao et al. 2011; Soltis et al. 2015).

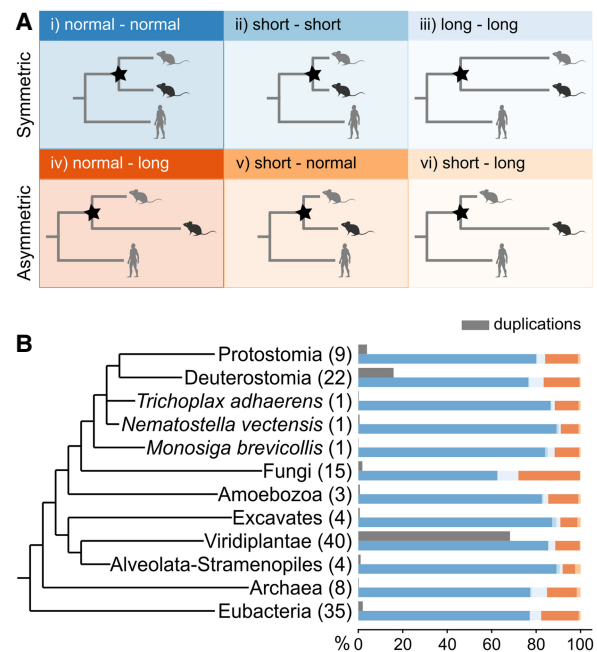
Further analysis of the distribution of duplications across the six categories within each major lineage (Fig. 2B, colored bars) revealed that the vast majority of duplications fell into the normal-normal group (~82%). The second most abundant category was the normal-long group (~13%), which includes paralogs with significantly different branch lengths. The remaining four categories (short-short, long-long, short-normal, and short-long) accounted for only ~6% of duplications. When grouping these categories based on their evolutionary pattern, 86% displayed symmetric evolution (both branches following similar rates: normal-normal, short-short, and long-long), with 14% exhibiting asymmetric evolution (one significantly longer: normal-long, short-normal, short-long). These results indicate that following duplication events, the

majority of paralogous pairs have similar evolutionary rates despite observed variations in absolute branch lengths.

Gene duplications in internal branches represent more ancient duplications compared to species-specific duplications (terminal branches). To investigate any differences, the distribution of duplications between these two types of duplications was examined within each major lineage subtree, excluding three major lineages which consist of only one species (*Trichoplax adhaerens*, *Nematostella vectensis*, *Monosiga brevicollis*). Across all cases, a large proportion of duplication events were present in terminal branches (~70%) (Supplemental Fig. S3A), which was similar in all major lineages (Supplemental Fig. S3B). When analyzing the distribution of the duplications between the two evolutionary models, a larger proportion of duplications followed symmetric evolution (Supplemental Fig. S3C) in all major lineages (Supplemental Fig. S3D). Furthermore, the proportion of duplications following asymmetric evolution was higher in internal than in terminal branches, suggesting that ancestral duplications may have more time to evolve divergent gene copies.

### LDOs show greater conservation of protein structure

Because protein structure is associated with protein function (Orengo et al. 1999), structural similarity can be used to test for functional divergence. However, structural similarity can be influenced by sequence similarity, because structural predictions often



**Figure 2.** Branch length analysis following duplication events. (A) Depiction of the six categories indicating differences in branch length following a duplication event (node marked with a star). The two rows divide the symmetric and asymmetric evolution of paralogs following duplication. The light silhouette (gray mouse) indicates the least diverged ortholog (LDO), and the dark silhouette (black mouse) indicates the most diverged ortholog (MDO). (B) Species tree showing the major lineages analyzed. The numbers in parentheses indicate the number of species included from the lineage. Gray bars represent the percentage of total duplications contributed by each lineage. Colored bars indicate the proportion of duplications within each category (color scheme as in A), relative to the total number of duplications in that lineage. The absolute numbers of duplications are provided in Supplemental Table S3.

rely on sequence information. To investigate this relationship, structural similarity (measured by Foldseek LDDT) and sequence similarity (percentage of identity) were calculated between paralogs classified into symmetric and asymmetric categories. A stronger correlation between structural and sequence similarity was observed in symmetric duplications (Pearson correlation coefficient [PCC]=0.75), whereas asymmetric duplications showed a weaker correlation (PCC=0.55) (Supplemental Fig. S4). This substantial difference between the two indicates that different evolutionary constraints or mechanisms act following symmetric versus asymmetric duplication.

To test the least diverged ortholog conjecture, controlling for sequence divergence, structural similarity was computed between each paralogous gene copy and a corresponding co-orthologous outgroup gene. The outgroups og1 and og2 were selected to have similar levels of sequence divergence to the MDO and the LDO, respectively (see Methods). To balance the need for comparable sequence divergence with the availability of suitable cases, the absolute difference in sequence identity between MDO-og1 and LDO-og2 was required to be at most 20%.

In asymmetric duplications, the median structural similarity for MDO-og1 pairs (0.723, dark orange) remained consistently lower than that for LDO-og2 pairs (0.7662, light orange) (Fig. 3A). In contrast, differences were minor for symmetric duplications (MDO-og1: 0.8506; LDO-og2: 0.8557). These findings suggest that, although sequence similarity influences structural similarity, structural differences may persist beyond what sequence alone can explain, potentially reflecting functional divergence.

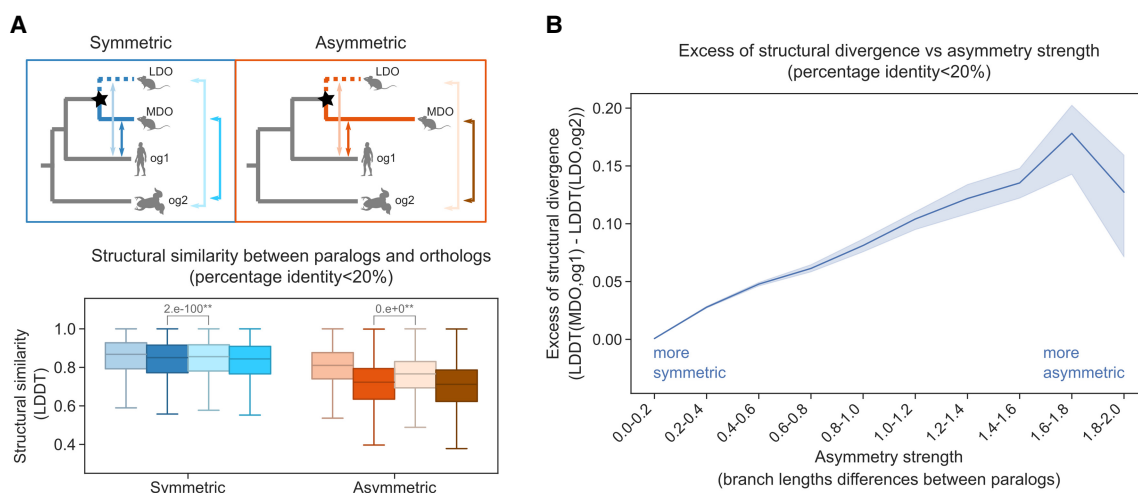
To further explore this, the asymmetry strength (measured as the observed branch length differences between paralogs) was compared to the excess structural divergence (calculated as  $LDDT[MDO, og1] - LDDT[LDO, og2]$ ) (Fig. 3B). The analysis revealed that greater asymmetry was associated with a stronger structural divergence of the MDO copy. A positive correlation was observed up to an asymmetry strength of approximately 1.8 (Fig. 3B). The decline in the excess of structural divergence may reflect the retention of conserved structural motifs despite the large sequence divergence.

To ensure that these results were not influenced by the choice of sequence identity difference threshold, additional analyses were performed using different distance similarity thresholds between 5% and 50%. In all cases, similar trends were observed (Supplemental Figs. S5, S6). Taken together, these results indicate that increased asymmetry in sequence evolution between paralogs is associated with greater structural divergence from the outgroup, even after controlling for sequence similarity. This supports the use of protein structure as a proxy for functional divergence. In line with the LDO conjecture, asymmetric duplications exhibit greater functional divergence, with the LDO copy more likely to retain the ancestral function.

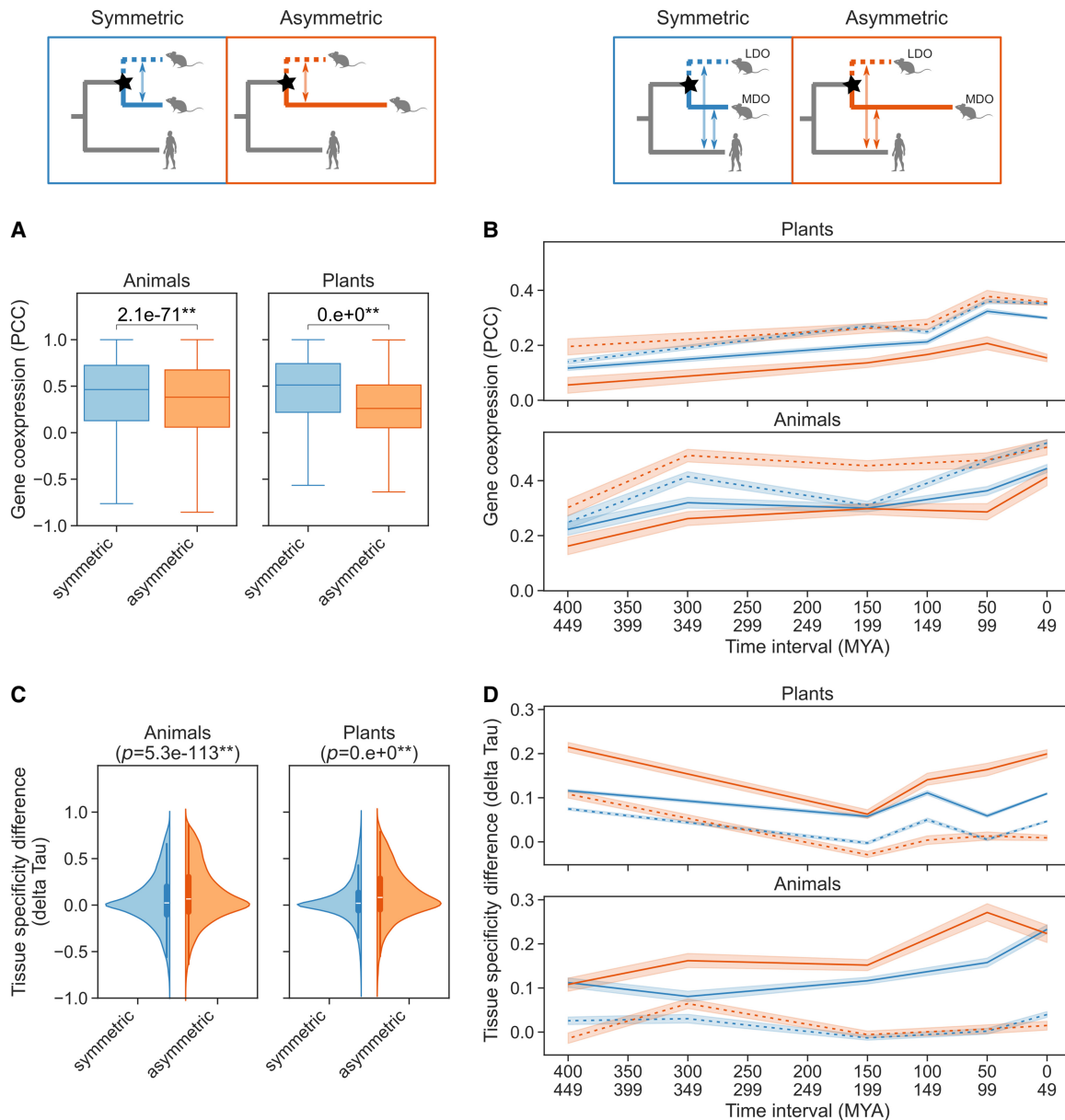
### Differences in branch lengths support differences in gene coexpression

To test functional differences across genes, gene coexpression was also used as a proxy for functional similarity. For animals, expression atlases of 16 species across 32 anatomical entities were obtained from the Bgee database (Bastian et al. 2021). For plants, publicly available expression atlases from 20 species across 72 anatomical samples were used (Supplemental Table S5; Julca et al. 2021; Koh et al. 2023). In order to assess if the observed differences in branch length following gene duplication have functional implications, gene coexpression between paralogs resulting from symmetric and asymmetric evolution was compared using the Pearson correlation coefficient. Subsequently, the direction of change was investigated by comparing the expression profiles of each paralog to their co-orthologous outgroup gene.

When analyzing pairs of paralogs in each category, within both animals and plants, paralogs from the asymmetric category showed lower correlation values compared to those from the symmetric category (Fig. 4A). Furthermore, because most duplications were located in terminal branches and represent more recent duplication events, their gene expression profiles were evaluated separately (Supplemental Fig. S7). In both cases, terminal and internal branches present consistent results, with paralogs in the symmetric category showing slightly higher coexpression



**Figure 3.** Structural similarity between paralogs and orthologs. (A) Structural similarity of each paralog to two different orthologs (og1 and og2) ( $n = 466,476$ ). Stars on the tree nodes indicate duplication events. Box colors correspond to the arrows in the figure, indicating the comparisons made. The  $P$ -value is indicated on top of the figure for the Mann-Whitney  $U$  test (\*\*\*)  $P < 0.01$ . All comparisons and  $P$ -values are indicated in Supplemental Table S4. (B) Comparison of the excess of structural divergence and the asymmetry strength ( $n = 466,476$ ). The shaded area represents the 95% confidence interval around the mean. The  $x$ -axis represents increasing bins of asymmetry strength.



**Figure 4.** Functional analysis of paralogs using expression data. (A) Gene coexpression values of paralogs following symmetric ( $n = 275,988$ ) and asymmetric ( $n = 31,671$ ) duplications, calculated using the Pearson correlation coefficient (PCC). The  $P$ -values are indicated on top of the figures for the Mann–Whitney  $U$  test. (\*)  $P < 0.05$ , (\*\*)  $P < 0.01$ . (B) Gene coexpression (PCC) between the LDO (dashed lines) and MDO (solid lines) genes and their closest out-group co-ortholog following symmetric (blue) and asymmetric (orange) duplications. The sample sizes for this plot are as follows: plants symmetric  $n = 114,702$ , plants asymmetric  $n = 13,533$ , animal symmetric  $n = 17,493$ , and animal asymmetric  $n = 7490$ . (C) Violin plot showing the tissue specificity differences (delta Tau) between the paralogs following symmetric (animals  $n = 77,889$ , plants  $n = 275,988$ ) and asymmetric (animals  $n = 17,205$ , plants  $n = 31,671$ ) duplications. (D) Tissue specificity differences (delta Tau) between the LDO (dashed lines) and MDO (solid lines) genes and their closest orthologs following symmetric (blue) and asymmetric (orange) duplications. The sample sizes are the same as in B. In plots B and D, the x-axis indicates the divergence time between the species compared in time intervals. Plant and animal data sets were compared independently. The comparisons used are shown above each panel, with duplication events marked by stars.

values. Notably, symmetrically evolving paralogs arising from duplications on the terminal branches within plants exhibited the highest coexpression (Supplemental Fig. S7). This observation may suggest that recent duplications have undergone less evolutionary divergence due to a shorter time frame. However, the same pattern is not observed in animals.

The results may be influenced by the different species' evolutionary history (e.g., polyploid plants). Therefore, each species

was analyzed independently and showed consistent results (Supplemental Fig. S8). Among the animals, 13 species showed similar patterns, indicating lower coexpression values between paralogs within the asymmetric category. However, three species, *Felis catus*, *Drosophila melanogaster*, and *Gorilla gorilla gorilla*, did not show significant differences between the two categories. For plants, 19 species supported the previous pattern, and only one (*Amborella trichopoda*) did not. Moreover, the polyploid plants

(*Brassica napus*, tetraploid; *Triticum aestivum*, hexaploid; *Glycine max*, tetraploid) exhibited major differences between symmetric and asymmetric paralogs. In summary, these results reveal a general trend where paralogs following asymmetric evolution exhibit less coexpression, suggesting that sequence divergence between paralogs correlates with differences in gene expression. Nevertheless, discrepancies observed in a few species highlight the importance of including many different species in this kind of study.

To investigate the direction of functional changes following duplication—whether the least or most diverged ortholog conserves the ancestral gene function—coexpression patterns were analyzed between each paralogous gene copy and their co-orthologous outgroup gene. For plants, 112,015 duplications were analyzed, with 99,995 symmetric and 12,020 asymmetric duplications. For animals, 22,885 duplications were examined, of which 16,273 were symmetric and 6612 were asymmetric duplications. Then, within each category, gene coexpression was computed (PCC) between each copy and the outgroup gene (Fig. 4B). Differences were observed in both symmetric and asymmetric categories, with greater differences observed between the asymmetric category (Fig. 4B). Moreover, in both data sets—animals and plants—the gene following the longest branch (MDO in the asymmetric category) consistently displays low coexpression values (Fig. 4B), suggestive of a lower functional similarity with the ancestral gene. Similar results were observed by analyzing each species independently (Supplemental Fig. S9).

To study the effect of time since duplication, coexpression was compared across species' divergence times (Fig. 4B). In both animals and plants, more ancient paralogs exhibited a slight decrease in coexpression.

In summary, these results suggest that paralogous genes descending from duplications with significantly different branch lengths have more distinct coexpression values. In particular, the most diverged orthologs display the greatest change in gene expression, thus suggesting that they have a greater functional divergence from the ancestral gene.

### MDOs show more functional specificity

Previous studies have shown that, after gene duplication, paralogs tend to increase the levels of tissue specificity, which occurs promptly after the duplication event (Huminięcki and Wolfe 2004; Huerta-Cepas et al. 2011). To test whether the least diverged ortholog (paralog with a shorter branch) or the most diverged ortholog (paralog with a longer branch) has different tissue specificity, the Tau score (defined in Methods) was computed for each gene in every species.

To compare each category, the difference in tissue specificity was calculated (delta Tau) between the most diverged ortholog and the least diverged ortholog in each group (symmetric and asymmetric duplications, Fig. 4C). A value of zero indicates no difference in sample-specificity between the compared genes, whereas positive values indicate that the MDO exhibits more specificity and negative values signify that the LDO is more specific. The results of this analysis show that paralogs within the asymmetric category, in both animals and plants, have slightly higher positive values, indicating that the MDOs tend to have a more tissue-specific pattern of expression (Fig. 4C). Similar results were observed when analyzing duplication events separately for internal and terminal branches, with terminal branches exhibiting more pronounced differences (Supplemental Fig. S10) and independently within each species (Supplemental Fig. S11). These results suggest

that the gene copy which exhibits greater sequence divergence tends to have greater tissue specificity.

To determine whether the most diverged ortholog also displays increased tissue specificity compared to the ancestral gene, the difference was computed between the Tau values for each paralogous gene copy and their co-orthologous outgroup gene (Fig. 4D). A zero value indicates no difference in tissue specificity, whereas positive values suggest increased specificity of the paralogous gene. In both categories, there are differences between the paralogs, with the MDO genes exhibiting higher delta Tau values. Specifically, in both plants and animals, the MDO gene of the asymmetric category consistently showed the greatest increase in tissue specificity over time (Fig. 4D). Moreover, similar results were observed when analyzing individual species (Supplemental Fig. S12). In addition, the delta Tau in animals showed a slight decrease over time, whereas in plants this decrease was not apparent (Fig. 4D). This result indicates that tissue specificity in animals may decrease over time, with recent duplications exhibiting higher specificity. To summarize, these findings suggest that when one gene copy is under selective pressure to retain the ancestral function (LDO), the other with the longer branch (MDO) tends to acquire a more specialized role, which may differ over time.

### Well-characterized gene families corroborate the LDO conjecture

To explore the findings of this study, well-characterized gene families from the literature with available experimental data were examined. The analysis revealed that symmetric duplications (normal-normal category) generally retain similar functions between paralogs, whereas asymmetric duplications (normal-long category) are more likely to diverge functionally. In these asymmetric pairs, the least diverged ortholog typically preserves the ancestral function.

A representative example is the *RNF113* gene family (Supplemental Fig. S13A), which duplicated independently in human and mouse. The human duplication is classified as normal-long, whereas the mouse duplication is normal-normal. In *Drosophila*, the ortholog *mdlc* is broadly expressed and essential for both neuronal development and spermatogenesis (Carney et al. 2013; Brattig-Correia et al. 2024). In humans, the LDO gene *RNF113A* retains broad expression and the ancestral neuronal role (Tsampoula et al. 2022), whereas *RNF113B* (MDO gene) is testis-specific and essential for meiosis (Brattig-Correia et al. 2024). In mouse, both *Rnf113a1* and *Rnf113a2* retain broad expression, yet *Rnf113a2* (MDO gene) is required for spermatogenesis and *Rnf113a1* (LDO gene) for neuronal development (Brattig-Correia et al. 2024), indicating some functional partitioning despite the classification as normal-normal.

Several clear examples of functional redundancy in normal-normal pairs were also identified. *SEC10a/b* in *Arabidopsis* (Supplemental Fig. S13B) share exocytosis-related functions (Vukašinović et al. 2014). The paralogs *Sar1a/b* in mouse (Supplemental Fig. S13C) are fully interchangeable for viability (Tang et al. 2024). For the paralogs *Cdx1/Cdx2* (Supplemental Fig. S13D), a knock-in of *Cdx2* rescues the *Cdx1*-null phenotype without defects (Savory et al. 2009).

In the normal-long category, *PRP18A/B* in *Arabidopsis* (Supplemental Fig. S13E) show functional divergence. *PRP18A* (LDO) retains a role in alternative splicing, similar to its yeast ortholog, whereas *PRP18B* (MDO) shows no detectable expression in 2-week-old seedlings, and T-DNA insertion mutants exhibit no observable phenotype compared to wild-type plants (Kanno et al.

2018). The paralogs *LAZY1/LAZY6* in *Arabidopsis* (Supplemental Fig. S13F) also illustrate this pattern: *LAZY1* (LDO) controls shoot orientation as in rice, whereas *LAZY6* (MDO) has a different expression profile and no significant phenotype in mutants (Li et al. 2007; Yoshihara and Spalding 2017; Waite and Dardick 2024). Similarly, studies in knockout mice suggest that *Becn1* and *Becn2* (Supplemental Fig. S13G) have distinct functions in regulating autophagosome formation and mitophagy (Galluzzi and Kroemer 2013; Quiles et al. 2023).

Overall, these examples are consistent with the broader patterns described in this study. Nevertheless, certain cases, such as the *Rnf113* family in mice, highlight the complexity of functional evolution postduplication and suggest that exceptions or intermediate scenarios exist. Whereas general trends across the Tree of Life were the focus of this study, gene-specific dynamics should be considered in future analyses.

## Discussion

To investigate the least diverged ortholog conjecture, a novel method was developed to identify paralogs with significantly divergent branch lengths, representing asymmetric evolution. This enabled the exploration of functional divergence between genes following duplication using gene trees from the PANTHER database, covering 143 species, along with structural data for more than one million proteins and expression atlases from 16 animal and 20 plant species. This study provides the most comprehensive analysis, using the largest data set to date for investigating functional divergence following gene duplication events. Furthermore, it avoids the pitfalls associated with comparing Gene Ontology functional annotations across different species (Altenhoff et al. 2012; Thomas et al. 2012; Gaudet and Dessimoz 2017).

In the PANTHER database, pairs of extant genes are classified as least diverged orthologs based on their gene tree branch lengths immediately following each postspeciation duplication event (Mi et al. 2010). However, this definition uses only the observed branch length, with no test being employed to determine the significance of any differences between the branch lengths following duplication. To address this gap, this study provides a generalization of the definition (least and most diverged orthologs), as well as a method to evaluate gene tree branch lengths and identify post-speciation duplication events with significant rate differences between branches. The results show that, despite observed differences in branch lengths between paralogs, most are not statistically significant (symmetric duplications). These results were consistent across all lineages included in this analysis. Previous studies have suggested that paralogs evolve at similar rates (Kondrashov et al. 2002), and only ~25% of duplications within a species display asymmetric evolution (Conant and Wagner 2003). Other studies have reported even lower occurrences, observing <20% exhibiting asymmetric evolution following duplication (Vance et al. 2022). A recent study observed that the most frequent fate of duplications is conservation (symmetry), followed by neofunctionalization (asymmetry) (Kalhor et al. 2024). The results presented here are consistent with previous studies, indicating that genes following a duplication event tend to evolve symmetrically. However, this observation does not preclude the possibility of functional changes occurring between paralogs.

If the least diverged ortholog conjecture holds true, it suggests that the least divergent copy (LDO) tends to retain the ancestral function, whereas the most divergent ortholog is free to evolve a

new function (Ohno 1970). Notably, in asymmetric duplications, the MDO copy exhibited more structural differences compared to the ancestral gene than the LDO. Even after accounting for sequence differences between paralogs and their orthologs, the MDO still showed an increase in structural divergence in asymmetric duplications. It has been observed that the function of a protein strongly depends on its structure (Orengo et al. 1999; Sotomayor-Vivas et al. 2022) and that an increase in structural similarity corresponds to an increase in functional similarity (Barrio-Hernandez et al. 2023). Additionally, whereas protein structure similarity can correlate with sequence similarity to some extent, it is well recognized that protein structures evolve more slowly than sequences. Moreover, structural changes may not always be proportional to sequence divergence due to compensatory mutations or structural constraints (Chothia and Lesk 1986; Gan et al. 2002; Illergård et al. 2009). Considering this, the results of this study suggest that sequence divergence, as reflected in gene tree branch lengths, correlates with functional divergence at the structural level.

Genes that coexpress are often members of similar biological processes and frequently share similar functions (Eisen et al. 1998; van Dam et al. 2018). In this context, coexpression analyses were used to test the LDO conjecture, showing that paralogous genes which have shorter branches (LDO) tend to coexpress with the orthologous gene, with larger differences observed in genes following asymmetric evolution.

Previous studies which focused on the models of gene duplication retention have proposed that either one (Gu et al. 2005; Brunet et al. 2006; Panchin et al. 2010; Pegueroles et al. 2013; Pich I Roselló and Kondrashov 2014) or both copies (Huminięcki and Wolfe 2004; Scannell and Wolfe 2008) undergo a rapid period of accelerated evolution, followed by a return to preduplication levels. Therefore, it has been suggested that expression divergence between paralogs can occur rapidly after a duplication event (Gu et al. 2002; Guschanski et al. 2017). Additionally, it has been observed that the expression of duplicated genes tends to evolve asymmetrically, with one copy maintaining the ancestral expression profile (Gu et al. 2005; Panchin et al. 2010; Pegueroles et al. 2013), and that asymmetric sequence divergence correlates with asymmetric functional divergence (Blanc and Wolfe 2004; Zhang et al. 2004; Kim and Yi 2006). Moreover, the rate of protein evolution may be influenced by expression levels, with highly expressed genes evolving more slowly due to stronger evolutionary constraints, such as protein misfolding (Drummond et al. 2005; Drummond and Wilke 2008). Overall, this study strongly corroborates previous works, supporting a rapid change in expression profile following duplication. Furthermore, the results highlight that the copy with the shortest branch (LDO), which evolves more slowly under major constraints, tends to retain the ancestral gene expression profile, further supporting the least diverged ortholog conjecture.

Assuming that the LDO conjecture holds and the MDO copy has a more divergent function, one hypothesis proposes that the new function tends to be more specific (Li et al. 2005). To test this “specialization hypothesis,” the tissue specificity of each gene copy was analyzed. This revealed that, following either symmetric or asymmetric evolution, the MDO copy tends to show higher tissue specificity, particularly notable in the MDO of asymmetric duplications. Additionally, for animals, the levels of tissue specificity in paralogs, when compared with their co-orthologs, appear to change over time, with the MDO tending to have a more specialized function in recent duplications. These results are consistent with previous studies, which showed that gene

duplication leads to a rapid increase in tissue specificity (Huminięcki and Wolfe 2004; Huerta-Cepas et al. 2011) and that expression divergence between paralogs results in increased tissue specificity (Assis and Bachtrog 2015). Furthermore, young paralogs were reported to be highly tissue-specific and become more broadly expressed with divergence time (Kryuchkova-Mostacci and Robinson-Rechavi 2016).

Genomic positional information has also been proposed as an additional indicator of functional retention after gene duplication. Specifically, genes that retain their ancestral genomic positions (“positional orthologs”) may be more likely to maintain their original function (Dewey 2011). This raises the possibility that, following duplication, the least diverged ortholog might not only be identified by sequence divergence but also by positional conservation. Although the present study did not incorporate positional information, considering both sequence and positional divergence could provide a more nuanced understanding of functional retention and represents an important direction for future work.

In conclusion, this comprehensive study using expression and structural aspects of function corroborates the least diverged ortholog conjecture: the gene copy with the shortest branch after a duplication event tends to retain the ancestral function. The LDO gene maintains a similar structure, expression profile, and tissue specificity to the ancestral gene, highlighting its functional conservation over time. Furthermore, the differences in structure and gene expression observed in paralogs following symmetric evolution suggested a rapid change occurring after gene duplication. Overall, these results significantly contribute to the understanding of gene duplication dynamics and their effect on gene function.

## Methods

### Gene tree processing

To study the least diverged ortholog conjecture, the species tree for 143 organisms and 15,693 gene trees were downloaded from PANTHER v18.0 (Thomas et al. 2022). PANTHER gene trees are annotated with evolutionary information in their internal nodes, including speciation, gene duplication, horizontal gene transfer (HGT), or unknown events (nodes for which the type cannot be reliably inferred). Specifically in bacteria, HGT events appear to be underestimated in this data set due to the methodology and the type of data used (single reference genome per bacterial species rather than a pangenome). Some specific branches in the gene trees were excluded from the analysis. These include branches with poorly fitting distances (arbitrarily set to 2.0 substitutions/site in PANTHER) and some duplication events that involve: duplications at the root of the gene tree (no speciation event to compare); those with horizontal gene transfer or “unknown” events annotated on either side of the duplication; as well as those with gene losses immediately following the duplication. Finally, a total of 15,577 gene trees were used for further analysis (Supplemental Table S2).

### Calculation of the expected gene tree branch lengths

Expected gene tree branch lengths were calculated for 2,648,252 branches across 15,577 gene trees using a simple evolutionary model: the expected length between two speciation events in a gene tree is equal to the corresponding branch in the species tree scaled by a gene family-specific rate. In the case of gene loss, the model assumes that the branch length in the gene tree represents

the sum of all the branches corresponding to the branch in the species tree scaled by the family-specific rate (Fig. 1A).

By expressing this as a system of approximate linear equations, it is possible to optimize the family rates and species tree branch lengths, constrained by the observed branch lengths within the many gene families. That is,  $|X| \approx \sum_{b_i \in X} |b_i|/f_R$ , where  $X$  represents a specific branch in a gene tree ( $|X|$  being the length of that branch), which corresponds to a set of branches in the species tree, each denoted  $b_i$  (with length  $|b_i|$ ), and  $f_R$  representing the family-specific factor.

The species tree used in PANTHER only provides topological relationships. As such, both the species tree branch lengths and family-specific factors were estimated from the observations of gene tree branch lengths. The method consists of two main steps:

1. *Species Tree Branch Fitting*: An iterative approach was used to fit the branch lengths of the species tree from the system of equations, while also optimizing the family-specific factors. Initially, species tree branches were set to the median of the lengths of matching branches from all the gene trees. This way, the species tree branch lengths reflect the typical amount of genetic change observed across all genes.

Using the initial species tree branch lengths, as described above, the family-specific factors were initialized according to the equations. The iterative procedure then has two steps: (1) refining the species tree branch lengths, using a least-squares solution, by estimating them while fixing the family-specific factor; (2) calculating the family-specific factors, while keeping the species tree branch lengths fixed. While optimizing the species tree branches, an additional normalization constraint was included, so that the average family-specific factor is 1. Convergence was achieved when the cosine distance between consecutive estimates is less than  $10^{-9}$ . The final estimate of species tree branch lengths was rounded to three decimal places. The species tree, including estimated branch lengths, is available in both Newick and PNG formats in the results folder of the git repository on GitHub ([https://github.com/DessimozLab/ldo\\_study/tree/main/results](https://github.com/DessimozLab/ldo_study/tree/main/results)) and the Supplemental Scripts.

2. *Family-Specific Factors*: Using the final estimate of species tree branch lengths, family-specific rates were calculated, without the additional normalization constraint used in step 1. The resulting family-specific factors are available in Supplemental Table S1.

In this way, the family-specific factor represents the best-fitting scaling factor that aligns the branch lengths of a gene tree with those of the species tree. It captures the overall rate of molecular evolution for each gene family, reflecting both neutral substitutions at tolerant sites and the effects of purifying selection at constrained sites. A family-specific factor greater than one indicates that the gene family is evolving faster—suggesting fewer functional constraints—whereas a factor less than one indicates slower evolution, consistent with stronger functional constraints compared to the average across families.

Using these species tree branch lengths and family-specific factors, expected branch lengths were calculated in the gene trees, before identifying statistically significant rate shifts (Supplemental Table S2). The Python scripts generated for this method are available on GitHub ([https://github.com/DessimozLab/ldo\\_study](https://github.com/DessimozLab/ldo_study)) and in the Supplemental Scripts.

Using these species tree branch lengths and family-specific factors, expected branch lengths were calculated in the gene trees, before identifying statistically significant rate shifts (Supplemental Table S2). The Python scripts generated for this method are available on GitHub ([https://github.com/DessimozLab/ldo\\_study](https://github.com/DessimozLab/ldo_study)) and in the Supplemental Scripts.

### Identification of branch length differences after duplication

Because all paralogs resulting from a duplication event can be classified as either least diverged orthologs (shorter branch) or most

diverged orthologs (longer branch), a statistical method was developed to determine whether the differences observed in the branch length postduplication were significant.

For 15,577 gene trees, the difference between observed and expected branch lengths (calculated under the model described in the previous section, from speciation event to speciation event) was computed (Fig. 1A). Then, to assess the statistical significance of the observed excess or depletion, these differences were converted into standardized Z-scores using the z-score function from the SciPy Python package v1.13.0 (Virtanen et al. 2020). An excess substitution can be interpreted as a potential signal of branch-specific adaptive selection and/or relaxed constraints. The resulting distributions of the normalized differences between branches with (876,261) and without (1,771,991) duplication were compared using the Mann–Whitney *U* test (also using SciPy). All further analysis filtered out gene trees without duplication events, which resulted in 13,597 gene trees and 364,829 duplication events distributed across different lineages of the Tree of Life.

For each duplication event, the two branches (Fig. 1A) were compared using their Z-score of the difference between observed and expected (from speciation to speciation event) by computing *P*-values and using a two-tailed test (with significance level  $\alpha=0.05$ ) to identify significantly shorter and longer than expected branches. When more than two gene copies were present, indicating either successive duplications, multifurcation (i.e., duplication nodes with more than two descendant branches), or both, the analysis was performed at every node starting from the leaves (Supplemental Fig. S14). Then, the least diverged ortholog was identified and compared to the other paralogs. This approach ensures that a long branch is included only once in the analysis.

Based on these comparisons, duplication events were placed into six categories (Fig. 2A): (i) normal-normal: both branches are not significantly different; (ii) short-short: both branches are significantly shorter than expected; (iii) long-long: both branches are significantly longer than expected; (iv) normal-long: only one branch is significantly longer than expected; (v) short-normal: only one branch is significantly shorter than expected; and (vi) short-long: one branch is significantly shorter than expected and the other is significantly longer than expected. Using this information, the data set was divided into two duplication models: symmetric (normal-normal, short-short, long-long) and asymmetric evolution (normal-long, short-normal, short-long). All subsequent analyses were based on this classification into these two categories: symmetric and asymmetric duplication events.

### Pairwise comparisons of structural and sequence identities

Structural data for 1,829,120 proteins out of 1,968,858 of the PANTHER data set were downloaded from the AlphaFold Protein Structure Database v4 (Varadi et al. 2022). The 139,738 remaining proteins did not have structural predictions available. Corresponding amino acid sequences for all proteins with available structural data were obtained from the UniProt database using the REST API (downloaded on 19 March 2025).

Pairwise structural similarity was computed for 926,814 paralog pairs using Foldseek version 8.ef4e960 (van Kempen et al. 2024). The Foldseek LDDT (Local Distance Difference Test) score reports the average LDDT of the alignment. LDDT scores express the percentage of interatomic distances and range from 0 (not conserved distances) to 1 (perfect model) (Mariani et al. 2013).

To estimate sequence similarity, multiple sequence alignments were generated for each gene family using MAFFT v7.526 (Katoh and Standley 2013). From these alignments, pairwise per-

centage identity between paralogs was calculated. The Pearson correlation coefficient was then computed to assess the relationship between sequence and structural similarity, separately for symmetric and asymmetric duplication categories.

In addition, pairwise sequence and structural similarity were calculated for each paralog and their co-orthologous outgroup gene (857,663 comparisons), using the approach described above. Two outgroup genes were selected to account for differences in the absolute branch length of each paralog. For the first outgroup (og1), co-orthologous genes were selected from the closest outgroup species before the duplication event in the gene tree. Then, the closest gene in that species (least diverged copy) was selected.

For the selection of the second outgroup (og2), the aim was to minimize the differences between LDO-og2 and MDO-og1, such that the sequence differences between MDO and og1 would be comparable to the differences between LDO and og2. The branch distance between MDO and og1 was calculated, and then an external ortholog with a similar distance to the LDO copy was selected as og2. To also account for sequence divergence, pairwise percentage identities between MDO and og1 and LDO and og2 were calculated using the multiple sequence alignment generated before (Supplemental Table S6).

For the analysis using the second outgroup (og2), the ideal scenario is that the difference in sequence similarity between MDO-og1 and LDO-og2 would be zero. In such cases, the comparison between MDO-og1 and LDO-og2 would be perfectly balanced in terms of sequence or branch distance. However, applying this strict filter drastically reduced the data set size (from 548,818 to 1319 (0.2%)). To address this, six more relaxed thresholds were used to filter for differences in percentage identity of: <5%, <10%, <20%, <30%, <40%, and <50%.

Distributions of the different data sets were compared using the Mann–Whitney *U* test.

### Asymmetry strength and excess of structural divergence analysis

To test whether the observed branch length differences between paralogs (referred to as asymmetry strength) correlate with structural divergence, the branch length differences for both symmetric and asymmetric duplications were calculated. Then, the differences in structural similarity (LDDT) between the comparisons LDO-og2 and MDO-og1 were computed (Supplemental Table S7). This analysis was performed for all data sets generated using the different sequence similarity and branch length thresholds described in the previous section.

### Compilation of gene expression atlases

As another functional characteristic of paralogous genes, expression atlases of 16 animals and 20 plants were used. For animals, RNA-seq gene expression data and anatomical annotations (UBERON terms) were retrieved from Bgee 15.1 (Bastian et al. 2021). This consisted of 2639 RNA-seq experiments across 32 different anatomical samples. For plants, expression atlases were taken from previous studies (The Plant Expression Omnibus) (Julca et al. 2021; Koh et al. 2023), and anatomical entities were obtained from the sample names of RNA-seq experiments in the NCBI BioSample database. Manual curation was conducted to define anatomical entity names (e.g., leaf, flowers), resulting in 72 different anatomical samples for plants. All species included, for both animals and plants, had at least four different anatomical samples (Supplemental Table S5). Also, only genes with TPM  $\geq 2$  in at least one sample were retained (plants = 446,852 and animals = 167,106 genes). Then, the TPM data of each gene was log-transformed

using the hyperbolic arcsine function (Johnson and Krishnan 2022):  $\arcsin h(x) = \ln(x + \sqrt{x^2 + 1})$ . Finally, the average value was computed for samples coming from the same anatomical entity.

### Gene coexpression analysis

To test for functional differences between pairs of paralogs (animals 95,094, and plants 307,659 pairs), as well as each paralog and the closest orthologous outgroup (animals 24,983, and plants 128,235 pairs), gene-gene coexpression patterns were computed using the Pearson correlation coefficient. To facilitate comparison across different species, only equivalent anatomical entities were used (e.g., leaf for *Arabidopsis* and tomato) (Supplemental Table S5). When selecting an orthologous outgroup to compare the paralogous genes to, the closest species was used for which data was available (Supplemental Table S8). Additionally, species were removed from the outgroup analysis when there were too few orthologous genes (<50 genes). Finally, distributions of the different data sets were compared using the Mann–Whitney *U* test.

### Analysis of gene expression specificity

Sample-specificity of genes, based on expression data, was calculated using the Tau score (Yanai et al. 2005; Kryuchkova-Mostacci and Robinson-Rechavi 2017). Tau values range from 0 to 1, where 0 indicates a gene is broadly expressed and 1 denotes tissue-specific expression. Gene pairwise comparisons were performed by calculating the difference in tissue specificity (delta Tau) between the two genes. Comparisons were performed for the pairs of extant paralogous genes (animals 95,094, and plants 307,659 pairs), as well as independently for each gene copy with the chosen orthologous outgroup gene (animals 24,983, and plants 128,235 pairs). Differences in their distribution were assessed using the Mann–Whitney *U* test.

### Estimation of divergence times

Divergence times (adjusted times) between species were obtained from the TimeTree database (Kumar et al. 2022) and are expressed in millions of years ago (MYA). These time estimates were used to plot the differences between genes in the expression profile and sample specificity analyses.

### Competing interest statement

The authors declare no competing interests.

### Acknowledgments

This work was funded by Swiss National Science Foundation grant 205085. I.J. acknowledges support by a Young Investigator Grant from the Novartis Foundation for Medical-Biological Research (24C173). We thank all members of the Comparative Genomics lab and Riccardo Delli Ponti for their valuable insights. We also thank the reviewers for their positive feedback.

*Author contributions:* C.D. and P.D.T. conceived the project idea. C.D., A.W.V., N.G., and I.J. designed the methodology. A.W.V. developed the scripts to analyze the data. A.W.V. and I.J. analyzed the data. I.J. and A.W.V. wrote the manuscript with comments from all authors. All authors read and approved the final manuscript.

### References

- Adams KL, Wendel JF. 2005. Novel patterns of gene expression in polyploid plants. *Trends Genet* **21**: 539–543. doi:10.1016/j.tig.2005.07.009
- Altenhoff AM, Studer RA, Robinson-Rechavi M, Dessimoz C. 2012. Resolving the ortholog conjecture: orthologs tend to be weakly, but significantly, more similar in function than paralogs. *PLoS Comput Biol* **8**: e1002514. doi:10.1371/journal.pcbi.1002514
- Assis R, Bachtrog D. 2013. Neofunctionalization of young duplicate genes in *Drosophila*. *Proc Natl Acad Sci* **110**: 17409–17414. doi:10.1073/pnas.1313759110
- Assis R, Bachtrog D. 2015. Rapid divergence and diversification of mammalian duplicate gene functions. *BMC Evol Biol* **15**: 138. doi:10.1186/s12862-015-0426-x
- Bandyopadhyay S, Sharan R, Ideker T. 2006. Systematic identification of functional orthologs based on protein network comparison. *Genome Res* **16**: 428–435. doi:10.1101/gr.4526006
- Barrio-Hernandez I, Yeo J, Jännes J, Mirdita M, Gilchrist CLM, Wein T, Varadi M, Velankar S, Beltrao P, Steinegger M. 2023. Clustering predicted structures at the scale of the known protein universe. *Nature* **622**: 637–645. doi:10.1038/s41586-023-06510-w
- Bastian FB, Roux J, Niknejad A, Comte A, Fonseca Costa SS, de Farias TM, Moretti S, Parmentier G, de Laval VR, Rosikiewicz M, et al. 2021. The Bgee suite: integrated curated expression atlas and comparative transcriptomics in animals. *Nucleic Acids Res* **49**: D831–D847. doi:10.1093/nar/gkaa793
- Blanc G, Wolfe KH. 2004. Functional divergence of duplicated genes formed by polyploidy during *Arabidopsis* evolution. *Plant Cell* **16**: 1679–1691. doi:10.1105/tpc.021410
- Brattig-Correia R, Almeida JM, Wyrwoll MJ, Julca I, Sobral D, Misra CS, Di Persio S, Guilgur LG, Schuppe H-C, Silva N, et al. 2024. The conserved genetic program of male germ cells uncovers ancient regulators of human spermatogenesis. *eLife* **13**: RP95774. doi:10.7554/eLife.95774
- Brunet FG, Roest Crollius H, Paris M, Aury J-M, Gibert P, Jaillon O, Laudet V, Robinson-Rechavi M. 2006. Gene loss and evolutionary rates following whole-genome duplication in teleost fishes. *Mol Biol Evol* **23**: 1808–1816. doi:10.1093/molbev/msl049
- Carney TD, Struck AJ, Doe CQ. 2013. *Midlife crisis* encodes a conserved zinc-finger protein required to maintain neuronal differentiation in *Drosophila*. *Development* **140**: 4155–4164. doi:10.1242/dev.093781
- Chothia C, Lesk AM. 1986. The relation between the divergence of sequence and structure in proteins. *EMBO J* **5**: 823–826. doi:10.1002/j.1460-2075.1986.tb04288.x
- Conant GC, Wagner A. 2003. Asymmetric sequence divergence of duplicate genes. *Genome Res* **13**: 2052–2058. doi:10.1101/gr.1252603
- Cusack BP, Wolfe KH. 2007. Not born equal: increased rate asymmetry in relocated and retrotransposed rodent gene duplicates. *Mol Biol Evol* **24**: 679–686. doi:10.1093/molbev/msl199
- Dewey CN. 2011. Positional orthology: putting genomic evolutionary relationships into context. *Brief Bioinformatics* **12**: 401–412. doi:10.1093/bib/bbr040
- Drummond DA, Wilke CO. 2008. Mistranslation-induced protein misfolding as a dominant constraint on coding-sequence evolution. *Cell* **134**: 341–352. doi:10.1016/j.cell.2008.05.042
- Drummond DA, Bloom JD, Adami C, Wilke CO, Arnold FH. 2005. Why highly expressed proteins evolve slowly. *Proc Natl Acad Sci* **102**: 14338–14343. doi:10.1073/pnas.0504070102
- Eisen MB, Spellman PT, Brown PO, Botstein D. 1998. Cluster analysis and display of genome-wide expression patterns. *Proc Natl Acad Sci* **95**: 14863–14868. doi:10.1073/pnas.95.25.14863
- Fitch WM. 1970. Distinguishing homologous from analogous proteins. *Syst Zool* **19**: 99–113. doi:10.2307/2412448
- Force A, Lynch M, Pickett FB, Amores A, Yan YL, Postlethwait J. 1999. Preservation of duplicate genes by complementary, degenerative mutations. *Genetics* **151**: 1531–1545. doi:10.1093/genetics/151.4.1531
- Gabalón T, Koonin EV. 2013. Functional and evolutionary implications of gene orthology. *Nat Rev Genet* **14**: 360–366. doi:10.1038/nrg3456
- Galluzzi L, Kroemer G. 2013. Common and divergent functions of Beclin 1 and Beclin 2. *Cell Res* **23**: 1341–1342. doi:10.1038/cr.2013.129
- Gan HH, Perlow RA, Roy S, Ko J, Wu M, Huang J, Yan S, Nicoletta A, Vafai J, Sun D, et al. 2002. Analysis of protein sequence/structure similarity relationships. *Biophys J* **83**: 2781–2791. doi:10.1016/S0006-3495(02)75287-9
- Gaudet P, Dessimoz C. 2017. Gene Ontology: pitfalls, biases, and remedies. *Methods Mol Biol* **1446**: 189–205. doi:10.1007/978-1-4939-3743-1\_14
- Gharib WH, Robinson-Rechavi M. 2013. The branch-site test of positive selection is surprisingly robust but lacks power under synonymous substitution saturation and variation in GC. *Mol Biol Evol* **30**: 1675–1686. doi:10.1093/molbev/mst062

- Gu Z, Nicolae D, Lu HH-S, Li WH. 2002. Rapid divergence in expression between duplicate genes inferred from microarray data. *Trends Genet* **18**: 609–613. doi:10.1016/S0168-9525(02)02837-8
- Gu X, Zhang Z, Huang W. 2005. Rapid evolution of expression and regulatory divergences after yeast gene duplication. *Proc Natl Acad Sci* **102**: 707–712. doi:10.1073/pnas.0409186102
- Guschanski K, Warnefors M, Kaessmann H. 2017. The evolution of duplicate gene expression in mammalian organs. *Genome Res* **27**: 1461–1474. doi:10.1101/gr.215566.116
- Huerta-Cepas J, Dopazo J, Huynen MA, Gabaldón T. 2011. Evidence for short-time divergence and long-time conservation of tissue-specific expression after gene duplication. *Brief Bioinformatics* **12**: 442–448. doi:10.1093/bib/bbr022
- Huminięcki L, Wolfe KH. 2004. Divergence of spatial gene expression profiles following species-specific gene duplications in human and mouse. *Genome Res* **14**: 1870–1879. doi:10.1101/gr.2705204
- Hurles M. 2004. Gene duplication: the genomic trade in spare parts. *PLoS Biol* **2**: E206. doi:10.1371/journal.pbio.0020206
- Illergård K, Ardell DH, Elofsson A. 2009. Structure is three to ten times more conserved than sequence—a study of structural response in protein cores. *Proteins* **77**: 499–508. doi:10.1002/prot.22458
- Jiao Y, Wickett NJ, Ayyampalayam S, Chanderbali AS, Landherr L, Ralph PE, Tomsho LP, Hu Y, Liang H, Soltis PS, et al. 2011. Ancestral polyploidy in seed plants and angiosperms. *Nature* **473**: 97–100. doi:10.1038/nature09916
- Johnson KA, Krishnan A. 2022. Robust normalization and transformation techniques for constructing gene coexpression networks from RNA-seq data. *Genome Biol* **23**: 1. doi:10.1186/s13059-021-02568-9
- Julca I, Ferrari C, Flores-Tornero M, Proost S, Lindner A-C, Hackenberg D, Steinbachová L, Michaelidis C, Gomes Pereira S, Misra CS, et al. 2021. Comparative transcriptomic analysis reveals conserved programmes underpinning organogenesis and reproduction in land plants. *Nat Plants* **7**: 1143–1159. doi:10.1038/s41477-021-00958-2
- Kalhor R, Beslon G, Lafond M, Scornavacca C. 2024. A rigorous framework to classify the postduplication fate of paralogous genes. *J Comput Biol* **31**: 815–833. doi:10.1089/cmb.2023.0331
- Kanno T, Lin W-D, Chang C-L, Matzke M, Matzke AJM. 2018. A genetic screen identifies PRP18a, a putative second step splicing factor important for alternative splicing and a normal phenotype in *Arabidopsis thaliana*. *G3 (Bethesda)* **8**: 1367–1377. doi:10.1534/g3.118.200022
- Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol* **30**: 772–780. doi:10.1093/molbev/mst010
- Kim S-H, Yi SV. 2006. Correlated asymmetry of sequence and functional divergence between duplicate proteins of *Saccharomyces cerevisiae*. *Mol Biol Evol* **23**: 1068–1075. doi:10.1093/molbev/msj115
- Koh E, Goh W, Julca I, Villanueva E, Mutwil M. 2023. PEO: Plant Expression Omnibus - a comparative transcriptomic database for 103 Archaeplastida. *Plant J* **117**: 1592–1603. doi:10.1111/tpj.16566
- Kondrashov FA, Rogozin IB, Wolf YI, Koonin EV. 2002. Selection in the evolution of gene duplications. *Genome Biol* **3**: RESEARCH0008. doi:10.1186/gb-2002-3-2-research0008
- Kryuchkova-Mostacci N, Robinson-Rechavi M. 2016. Tissue-specificity of gene expression diverges slowly between orthologs, and rapidly between paralogs. *PLoS Comput Biol* **12**: e1005274. doi:10.1371/journal.pcbi.1005274
- Kryuchkova-Mostacci N, Robinson-Rechavi M. 2017. A benchmark of gene expression tissue-specificity metrics. *Brief Bioinformatics* **18**: 205–214. doi:10.1093/bib/bbw008
- Kumar S, Suleski M, Craig JM, Kaspric AE, Sanderford M, Li M, Stecher G, Hedges SB. 2022. TimeTree 5: an expanded resource for species divergence times. *Mol Biol Evol* **39**: msac174. doi:10.1093/molbev/msac174
- Leitch AR, Leitch IJ. 2008. Genomic plasticity and the diversity of polyploid plants. *Science* **320**: 481–483. doi:10.1126/science.1153585
- Li W-H, Yang J, Gu X. 2005. Expression divergence between duplicate genes. *Trends Genet* **21**: 602–607. doi:10.1016/j.tig.2005.08.006
- Li P, Wang Y, Qian Q, Fu Z, Wang M, Zeng D, Li B, Wang X, Li J. 2007. LAZY1 controls rice shoot gravitropism through regulating polar auxin transport. *Cell Res* **17**: 402–410. doi:10.1038/cr.2007.38
- Lynch M, Force A. 2000. The probability of duplicate gene preservation by subfunctionalization. *Genetics* **154**: 459–473. doi:10.1093/genetics/154.1.459
- Mariani V, Biasini M, Barbato A, Schwede T. 2013. IDDT: a local superposition-free score for comparing protein structures and models using distance difference tests. *Bioinformatics* **29**: 2722–2728. doi:10.1093/bioinformatics/btt473
- Mi H, Dong Q, Muruganujan A, Gaudet P, Lewis S, Thomas PD. 2010. PANTHER version 7: improved phylogenetic trees, orthologs and collaboration with the Gene Ontology Consortium. *Nucleic Acids Res* **38**: D204–D210. doi:10.1093/nar/gkp1019
- Nehrt NL, Clark WT, Radivojac P, Hahn MW. 2011. Testing the ortholog conjecture with comparative functional genomic data from mammals. *PLoS Comput Biol* **7**: e1002073. doi:10.1371/journal.pcbi.1002073
- Ohno S. 1970. *Evolution by gene duplication*. Springer, Berlin, Heidelberg.
- Ohno S. 1972. An argument for the genetic simplicity of man and other mammals. *J Hum Evol* **1**: 651–662. doi:10.1016/0047-2484(72)90011-5
- Orongo CA, Todd AE, Thornton JM. 1999. From protein structure to function. *Curr Opin Struct Biol* **9**: 374–382. doi:10.1016/S0959-440X(99)80051-7
- Panchin AY, Gelfand MS, Ramensky VE, Artamonova II. 2010. Asymmetric and non-uniform evolution of recently duplicated human genes. *Biol Direct* **5**: 54. doi:10.1186/1745-6150-5-54
- Pegueroles C, Laurie S, Albà MM. 2013. Accelerated evolution after gene duplication: a time-dependent process affecting just one copy. *Mol Biol Evol* **30**: 1830–1842. doi:10.1093/molbev/mst083
- Pich I, Roselló O, Kondrashov FA. 2014. Long-term asymmetrical acceleration of protein evolution after gene duplication. *Genome Biol Evol* **6**: 1949–1955. doi:10.1093/gbe/evu159
- Quiles JM, Najor RH, Gonzalez E, Jeung M, Liang W, Burbach SM, Zumaya EA, Diao RY, Lampert MA, Gustafsson ÅB. 2023. Deciphering functional roles and interplay between Beclin1 and Beclin2 in autophagosome formation and mitophagy. *Sci Signal* **16**: eabo4457. doi:10.1126/scisignal.abo4457
- Savory JGA, Pilon N, Grainger S, Sylvestre J-R, Béland M, Houle M, Oh K, Lohnes D. 2009. Cdx1 and Cdx2 are functionally equivalent in vertebral patterning. *Dev Biol* **330**: 114–122. doi:10.1016/j.ydbio.2009.03.016
- Scannell DR, Wolfe KH. 2008. A burst of protein sequence evolution and a prolonged period of asymmetric evolution follow gene duplication in yeast. *Genome Res* **18**: 137–147. doi:10.1101/gr.6341207
- Sehrish T, Symonds VV, Soltis DE, Soltis PS, Tate JA. 2014. Gene silencing via DNA methylation in naturally occurring *Tragopogon miscellus* (Asteraceae) allopolyploids. *BMC Genomics* **15**: 701. doi:10.1186/1471-2164-15-701
- Sémon P, Wolfe KH. 2007. Consequences of genome duplication. *Curr Opin Genet Dev* **17**: 505–512. doi:10.1016/j.gde.2007.09.007
- Smith JM, Smith NH. 1996. Synonymous nucleotide divergence: what is “saturation”? *Genetics* **142**: 1033–1036. doi:10.1093/genetics/142.3.1033
- Soltis PS, Marchant DB, Van de Peer Y, Soltis DE. 2015. Polyploidy and genome evolution in plants. *Curr Opin Genet Dev* **35**: 119–125. doi:10.1016/j.gde.2015.11.003
- Sotomayor-Vivas C, Hernández-Lemus E, Dorantes-Gilardi R. 2022. Linking protein structural and functional change to mutation using amino acid networks. *PLoS One* **17**: e0261829. doi:10.1371/journal.pone.0261829
- Tang VT, Xiang J, Chen Z, McCormick J, Abbineni PS, Chen X-W, Hoenerhoff M, Emmer BT, Khoriaty R, Lin JD, et al. 2024. Functional overlap between the mammalian *Sar1a* and *Sar1b* paralogs in vivo. *Proc Natl Acad Sci* **121**: e2322164121. doi:10.1073/pnas.2322164121
- Thomas PD, Wood V, Mungall CJ, Lewis SE, Blake JA, Gene Ontology Consortium. 2012. On the use of Gene Ontology annotations to assess functional similarity among orthologs and paralogs: a short report. *PLoS Comput Biol* **8**: e1002386. doi:10.1371/journal.pcbi.1002386
- Thomas PD, Ebert D, Muruganujan A, Mushayama T, Albu L-P, Mi H. 2022. PANTHER: making genome-scale phylogenetics accessible to all. *Protein Sci* **31**: 8–22. doi:10.1002/pro.4218
- Tsampoula M, Tarampoulou I, Manolakou T, Ninou E, Politis PK. 2022. The neurodevelopmental disorders associated gene *Rnf113a* regulates survival and differentiation properties of neural stem cells. *Stem Cells* **40**: 678–690. doi:10.1093/stmcls/sxac030
- Vance Z, Niezabitowski L, Hurst LD, McLysaght A. 2022. Evidence from *Drosophila* supports higher duplicability of faster evolving genes. *Genome Biol Evol* **14**: evac003. doi:10.1093/gbe/evac003
- van Dam S, Vösa U, van der Graaf A, Franke L, de Magalhães JP. 2018. Gene co-expression analysis for functional classification and gene-disease predictions. *Brief Bioinformatics* **19**: 575–592. doi:10.1093/bib/bbw139
- van Kempen M, Kim SS, Tumescheit C, Mirdita M, Lee J, Gilchrist CLM, Söding J, Steinegger M. 2024. Fast and accurate protein structure search with Foldseek. *Nat Biotechnol* **42**: 243–246. doi:10.1038/s41587-023-01773-0
- Varadi M, Anyango S, Deshpande M, Nair S, Natassia C, Yordanova G, Yuan D, Stroe O, Wood G, Laydon A, et al. 2022. AlphaFold Protein Structure Database: massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic Acids Res* **50**: D439–D444. doi:10.1093/nar/gkab1061
- Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, Burovski E, Peterson P, Weckesser W, Bright J, et al. 2020. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat Methods* **17**: 261–272. doi:10.1038/s41592-019-0686-2
- Vukašinović N, Cvrčková F, Eliáš M, Cole R, Fowler JE, Žárský V, Synek L. 2014. Dissecting a hidden gene duplication: the *Arabidopsis thaliana* SEC10 locus. *PLoS One* **9**: e94077. doi:10.1371/journal.pone.0094077

Warwick Vesztröcy et al.

---

- Waite JM, Dardick C. 2024. IGT/LAZY genes are differentially influenced by light and required for light-induced change to organ angle. *BMC Biol* **22**: 8. doi:10.1186/s12915-024-01813-4
- Wang Y, Wang X, Paterson AH. 2012. Genome and gene duplications and gene expression divergence: a view from plants. *Ann N Y Acad Sci* **1256**: 1–14. doi:10.1111/j.1749-6632.2011.06384.x
- Yanai I, Benjamin H, Shmoish M, Chalifa-Caspi V, Shklar M, Ophir R, Bar-Even A, Horn-Saban S, Safran M, Domany E, et al. 2005. Genome-wide midrange transcription profiles reveal expression level relationships in human tissue specification. *Bioinformatics* **21**: 650–659. doi:10.1093/bioinformatics/bti042
- Yoshihara T, Spalding EP. 2017. LAZY genes mediate the effects of gravity on auxin gradients and plant architecture. *Plant Physiol* **175**: 959–969. doi:10.1104/pp.17.00942
- Zhang Z, Gu J, Gu X. 2004. How much expression divergence after yeast gene duplication could be explained by regulatory motif evolution? *Trends Genet* **20**: 403–407. doi:10.1016/j.tig.2004.07.006

Received October 29, 2024; accepted in revised form August 7, 2025.