

# **GageTracker: a tool for dating gene age by micro- and macro-synteny with high speed and accuracy**

Chengchi Fang<sup>1,†</sup>, Chuan Dong<sup>2,†</sup>, Cheng Wang<sup>1,3</sup>, Fan Xiong<sup>1</sup>, Suxiang Lu<sup>1</sup>, Wenyu Fang<sup>3</sup>, Tong Li<sup>1,3</sup>, Xiaoni Gan<sup>1</sup>, Liandong Yang<sup>1</sup>, Honghui Zeng<sup>1,\*</sup>, Shunping He<sup>1,4,\*</sup>

<sup>1</sup> The Key Laboratory of Aquatic Biodiversity and Conservation of Chinese Academy of Sciences, Institute of Hydrobiology, Chinese Academy of Sciences, Wuhan, Hubei, 430072, China

<sup>2</sup> State Key Laboratory of Subtropical Silviculture, Zhejiang A&F University, Hangzhou, Zhejiang, 311300, China

<sup>3</sup> University of Chinese Academy of Sciences, Beijing, 100039, China

<sup>4</sup> State Key Laboratory of Freshwater Ecology and Biotechnology, Institute of Hydrobiology, Chinese Academy of Sciences, Wuhan, Hubei, 430072, China

† The authors would like to make it clear that, in their view, the first two authors should be considered as CO-FIRST AUTHORS.

\* To whom correspondence should be addressed. Email: clad@ihb.ac.cn, and zhh@ihb.ac.cn.

## 24 ABSTRACT

25 With the advent of the Earth Genome Project, an increasing number of species' genomes presents  
 26 exciting opportunities for exploring genetic and phenotypic diversity in organisms. Determining the  
 27 origin time of genes facilitates the elucidation of crucial genetic mechanisms underlying significant  
 28 biological evolutionary questions such as the transition from aquatic to terrestrial life, the emergence  
 29 of mammals, the origin of humans, as well as the development of species- or lineage-specific traits.  
 30 However, accurately determining the origin time of these genes in species separated by long  
 31 evolutionary distances remains a major challenge in bioinformatics as these genes often undergo  
 32 significant changes in their genome sequences, making it difficult to trace them back to their origin.  
 33 Here, we proposed a new approach for dating gene age based on the micro- and macro-synteny  
 34 algorithms. This approach employs the parallel computation of orthologous genome alignments across  
 35 multiple species. Our method was integrated into the GageTracker (Gene Age Tracker) software,  
 36 providing a fast and accurate way to trace gene age with minimal user input, available at  
 37 <https://github.com/RiversDong/GageTracker>. Benchmarked against the simMammals dataset  
 38 (Alignathon), GageTracker achieved the same high-quality genome alignments as the optimized LastZ  
 39 aligner, but improved operation speed by 1.4-7 times. In a separate analysis of 12 *Drosophila*  
 40 genomes, GageTracker efficiently assessed the ages of 23,720 genes (including ~13,965 protein-  
 41 coding genes) in just ~22 hours with default parameters. When comparing with the GenTree database  
 42 (recognized as the most comprehensive and accurate tool for evaluating gene age), GageTracker  
 43 achieved an impressive ~94.4% accuracy and ~99% macro consistency in assessing the age of  
 44 protein-coding genes. Moreover, for the ~5.6% conflicting genes, GageTracker displayed slightly  
 45 higher support rates than GenTree, as evidenced by data from OrthoDB, FlyBase, and Ensembl  
 46 ortholog databases. Notably, younger genes identified by GageTracker exhibited a preferential  
 47 expression pattern in the testis, further reinforcing the reliability of GageTracker in accurately tracing  
 48 gene age.

49

50

51

## 52 INTRODUCTION

53 Since Darwin's era, biologists have been keenly focused on the question of how organisms evolved  
 54 from a common ancestor to the diverse species that we see today. Genomic variation has been a major  
 55 contributor to this vast biological diversity (1). With the advent of the Earth Genome Project,  
 56 including the Bird 10K Genomes, Fish 10K Genomes, and Protist 10K Genomes, the analysis of an  
 57 increasing number of species' genomes presents exciting opportunities for exploring genetic and  
 58 phenotypic diversity in organisms. Determining the origin time of genes facilitates the elucidation of  
 59 crucial genetic mechanisms underlying the development of species- or lineage-specific traits.  
 60 Researchers have discovered that even closely related species exhibit a certain degree of variability in  
 61 genome size and gene number (2). The abundance of available genomes provides a valuable  
 62 opportunity to answer several central questions in evolution, especially when it comes to  
 63 comprehending the evolution of new genes. Approximately 380 new functional genes related to brain  
 64 development evolved in a relatively short period after the divergence of humans and chimpanzees (~6  
 65 million years ago). These genes played a crucial role in enhancing human cognition and helping to  
 66 shed light on what makes humans as human being (3-5). The origin of novel genes has been  
 67 considered to be a major contributor to adaptive innovation (1,6-8). They often result in the  
 68 development of new adaptive traits and have a profound impact on resolving sexual conflicts,  
 69 embryonic development, and reproduction (9-11). In flies, the origin of new genes has provided a new  
 70 genetic basis to species and they are frequently highly expressed in reproductive organs, such as  
 71 *jingwei* and *sphinx* (11,12). The roles of new genes that could be involved in the important biological  
 72 process were also demonstrated by a series of high-throughput RNA silencing (13) and CRISPR-Cas9  
 73 experimental technology (14). In plants, new genes have been shown to play important roles in  
 74 various biological processes, such as male sterility and plant height (15,16). It was also found that  
 75 fusion genes could regulate phenotype traits including seed germination, shoot length, and root length  
 76 (17). Furthermore, Zhang et al. discovered that new protein-coding genes could emerge from non-  
 77 coding DNA sequences, providing a source for additional diversity in protein composition (18). The  
 78 new genes in bamboo were related to the origin and evolution of its unique character of rapid growth  
 79 (19). In prokaryotic species, the evolution of new genes, which was shaped by the GC content,  
 80 expression level, as well as the essentiality of a gene, were important factors that could influence the  
 81 evolutionary rate (20,21). Together, more and more studies have shown that new genes were one of

82 the main driving forces of evolution and innovation, and promote organisms to produce special  
83 phenotypic and physiological characteristics. Therefore, the annotation of new genes contributed  
84 greatly to the understanding of mechanisms underlying these important biological functions.

85 However, the progress in comprehending the impact of new genes on evolutionary traits and  
86 lineage-specific characteristics has been hindered by the absence of a rapid and precise approach to  
87 identifying orthologs and determining the origination age of genes. The annotation of gene age via  
88 comparative genomics has been primarily achieved through three strategies: 1) the homologous gene  
89 family-level dating strategy (22-24); 2) the synteny-based gene-level dating strategy (22,25-27); 3) the  
90 comprehensive orthologous database (4,24,28,29). The first strategy was significantly dependent on  
91 the quality of gene annotations in the reference species, thereby leading to potential false positives in  
92 the new gene catalog in cases where the annotation was absent in the reference species. The second  
93 strategy relied heavily on genome synteny analysis, which allowed for the avoidance of imperfect  
94 protein gene annotations and ensured high accuracy in detecting new genes. Nonetheless, this strategy  
95 had certain constraints in managing complex sequences such as genome sequencing gaps and low-  
96 complexity regions (e.g., telomeres), leading to potential inaccuracies. On the other hand, the third  
97 approach, encompassing GenTree, GenOrigin, and OrthoDB, typically employed a combination of  
98 methods to construct a searchable database for tracing gene age (4,24,30). However, these databases  
99 had limited coverage, given that they only included a few well-studied model species and could not be  
100 updated or expanded to include new species of interest, thereby impeding their utility, particularly in  
101 an era of rapid growth in genomic data.

102 The rapid accumulation of genomic data has created a challenge for gene age annotation.  
103 Thousands of genome assemblies are publicly available and additional species are being sequenced as  
104 part of projects, further exacerbating the challenge (31-33). This dense sampling and sequencing of  
105 the tree of life are set to provide a comprehensive understanding of evolution and biodiversity. With  
106 the rapid expansion of genomic data to fill out the phylogeny of vertebrates, there is going to be a  
107 strong desire to explore genetic innovation and phenotypic diversity of the different species,  
108 meanwhile, the growing data also provided a great opportunity for tracing gene origin time in more  
109 details. Thus, we proposed a new tool, GageTracker, for tracing gene age in the era of rapid expansion  
110 of genomic data. Based on the micro- and macro-synteny algorithm, GageTracker was a one-  
111 command running software to search ortholog genome alignments suitable for multiple species and

allow a fast and accurate trace gene age with minimal user inputs. It obtained a high alignment quality as the optimized LastZ software but significantly saved the running time as well. Compared with the GenTree database, it achieved ~94.4% accurate consistency and ~99% macro consistency in assessing the age of protein-coding genes. For the ~5.6% conflicting gene age, GageTracker also showed a slightly higher support rate from orthoDB, FlyBase, and Ensembl ortholog database than the Gentree database, suggesting its reliability for dating gene age.

## MATERIAL AND METHODS

### General Algorithm of GageTracker

The workflows of GageTracker were mainly divided into three parts: data inputs, genome alignment and processing, and age-infer (Figure 1A). The genome alignment behavior of GageTracker was specified by the configuration file and command line parameters. The Last v1229 program (34) was used to conduct genome alignment, and several toolkits from UCSC (mainly including maftochain, chainCleaner, chainSwap, netSyntenic, netChainSubset, chainStitchId, and chainNet) were employed to detect the Reciprocal Best Syntenic Net hits (RBH) alignments according to the description of RBH pipeline (doRecipBest.pl), which can be available at [https://genomewiki.ucsc.edu/index.php?title=HowTo:\\_Syntenic\\_Net\\_or\\_Reciprocal\\_Best](https://genomewiki.ucsc.edu/index.php?title=HowTo:_Syntenic_Net_or_Reciprocal_Best). We used the maf2synteny (35) and the bedtools (36) to construct the block regions and detect sequencing gaps (left panel of Figure 1B). Finally, GageTracker labels gene age according to a voting strategy (left and middle panel of Figure 1B).

Given a genome that has evolved through a steady-state process of gene acquisition and loss, we have employed a voting strategy to assess whether a gene originated on a particular branch (left and middle panel of Figure 1B). First, we conducted a gene presence analysis on each branch. Specifically, we inferred the presence of a gene within a branch if it had orthologous segments in any species within that branch. Subsequently, we systematically evaluated each evolutionary branch in chronological order, from the young to the old, and computed the support rate of the gene's origin on that branch. For example, Gene 1 is found to have orthologs on branches 6, 3, 2, and 0, indicating its presence on branches 6, 3, 2, and 0 but absence on branches 5, 4, and 1 (orange rectangle: presence,

blank rectangle: absence, left and middle panel of Figure 1B). Then, we will proceed by sequentially calculating the support rate on branches 6, 3, 2, and 0 (percentage in orange rectangles, left and middle panel of Figure 1B). When focusing on branch 6, Gene 1's support rate is 100%; when focusing on branch 3, its support rate is 50% (2/4, present in 2 branches (6, 3) and absent in 2 branches (5, 4)); when focusing on branch 2, its support rate is 60% (3/5, present in 3 branches (6, 3, 2) and absent in 2 branches (5, 4)); when focusing on branch 0, its support rate reaches 57% (4/7, present in 4 branches (6, 3, 2, 0), absent in 3 branches (5, 4, 1)). Therefore, Gene 1 is assigned to the oldest branch 0, as the support rate is above the voting threshold of 50%. Similarly, Gene 2 is assigned to the oldest branch 1. At the same time, Gene 3 appears on branches 6 and 0. But its support rate of originating on branch 0 (2/7, present in 2 branches, absent in other 5 branches) is only 28.5%, indicating that Gene 3 is more likely to have originated independently rather than being lost in all intermediate species. Thus, Gene 3 independently emerged on branch 6 as a newly originated gene.

In principle, this algorithm could well distinguish independent origin genes to avoid overestimating genes' age. The user could adjust the voting parameter in the control file (voting=0.5) according to the divergence time between species and their own needs.

### **The masked strategy for genomes**

Soft masking of repeat sequences was a common preprocessing step in genome alignment and gene age estimation. GageTracker utilized a *de novo* prediction method for masking genomic repetitive sequences. It utilized WindowMasker to mask interspersed repeats sequences and Tantan to simple-sequence from scratch in genome sequences. In GageTracker, the repeat sequences were first softly masked. Because the protein-coding sequence contained the conserved genetic information between genes. To prevent masking such functional regions, we converted the incorrectly masked sequence to uppercase in the protein-coding regions.

### **Benchmark genomes alignment**

To evaluate the performance of GageTracker, we included two types of benchmark genomes: simulated data and real data, respectively. The simulated mammal's data was downloaded from the packageMammals in Alignathon database (37). It contained five simulated benchmarks (simHuman, simMouse, simRat, simCow, simDog) that consisted of genomes of different species, sizes, and

complexities. The real data sets, including *Homo sapiens*, *Mus musculus*, and *Drosophila melanogaster*, as well as their GTF annotation files, were retrieved from Ensembl (version 108). Other 11 *Drosophila* genomes were obtained from FlyBase (version FB2022\_06). The detailed version number of each genome were listed in Supplementary Table S1.

The genome alignments were performed by GageTracker, Lastal, and PlastZ aligners respectively on the same computational platform using commands that listed in the following script in the calculation section. Each aligner was controlled in 6 threads on the servers with an E7-4850-CPU (2.5GHZ) array and 1 TB of memory.

- Only running genome alignments in GageTracker command line:  
*GageTracker example.ctl -step1*
- Lastal commands line, which contains two steps: convert repeat sequence into uppercase, and perform genome alignment by Lastal  

```
less mus.fasta|perl -e '{while(<>){if(/^>/){print;}else{$_=uc($_); print}}}'> mus.up.fasta
less human.fasta|perl -e '{while(<>){if(/^>/){print;}else{$_=uc($_); print}}}'> human.up.fasta
lastdb -c -P6 -R10 -uMAM8 humandb ./human.up.fasta
lastal -P6 -C2 -u0 -m50 -p HOXD70 ./humandb ./mus. up.fasta
```
- Perform genome alignment by Lastz  

```
python3.7 PLastZ.py mus.fasta Human.fasta ./output -lo="O=400 E=30 L=3000 H=2200 T=1 --ambiguous=iupac --format=lav" -p 6
```

## Evaluation of alignments

Genome alignment results often yield many-to-many matches due to the prevalence of gene duplication as a major driver of new gene birth. However, we focused specifically on identifying one-to-one orthologous alignments. To obtain our targeted regions, the many-to-many alignments generated by the aligners were fed into a similar doRecipBest.pl pipeline from UCSC, which mainly included mafToPsl, pslToChain, chainCleaner, chainStitchId, chainSwap, chainNet, netSyntenic, netChainSubset, chainStitchId, chainSwap, chainNet, and netToAxt toolkits. This could create the query-reference reciprocal best chains and nets files (refer to one-to-one orthologous alignments in this work).

Second, we processed the ancestor.noparalogies.maf data in the benchmark alignment file of the Alignathon database by mafTools pipeline (37). It mainly contained: i) The mafFilter to extract alignments that only match the query-reference species; and ii) The mafTransitiveClosure to perform the transitive closure on these query-reference alignments. This could produce four groups (simCow-simHuman, simDog-simHuman, simMouse-simHuman, simRat-simHuman) of query-reference

genome alignments and the simHuman was the reference genome. These four groups' genome alignments were considered one-to-one truth alignments in our work.

$$\text{Precision} = \frac{|A \cap T|}{|A|}, \text{Recall} = \frac{|A \cap T|}{|T|}, F\text{-score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (1)$$

$$\text{OP} = \frac{OG}{\text{Min}(TG \text{ in } \text{GagTracker}, TG \text{ in } \text{GenTree})} \quad (2)$$

Finally, we used the mafComparator to compare the test alignments to these benchmark alignments and calculated the number of identical base pairs. To assess the quality of the three aligners (Lastal, GageTracker, and LastZ), we calculated recall, precision, and balanced *F*-score values in the usual way as equation (1). Let *T* denote the number of base pairs in one-to-one truth alignments, and *A* denote the number of base pairs in one-to-one test alignments. Then the precision was the ratio of test positive predictions (*A*∩*T*) to the total number of test alignments (*A*); recall was the ratio of test positive predictions (*A*∩*T*) to the total number of truth alignments (*T*); and *F*-score was the harmonic mean of recall and precision. It is the trade-off between precision and recall and provides a single metric to evaluate the overall performance of a classification model.

### Age comparison in three different levels

We used the dating gene age of *D. melanogaster* as a case to compare the performance between GageTracker and GenTree at three different levels: the global level, the nearly-detailed level, and the detailed level. More specifically, the first level was used to measure the global consistency with GenTree (only focusing on the proportion of the genes with the same gene age). In the nearly-detailed level, we divided genes into two categories young genes (from branch 1 to branch 6) and old genes (genes in branch 0), and calculated the overlapping proportion (OP) of the two annotated age versions (GenTree and GageTracker). The overlapping proportion was calculated by equation (2), where OG represents the number of overlapping genes in the same comparison level and TG was the total genes in the current comparison level.

In addition, we estimated the macro-consistency between GageTracker and GenTree by comparing the branch distance of the same gene (branch distance = |Age in GageTracker – Age in GenTree|). There are seven branches in the evolutionary lineage in our case study, and half of them are three. We considered a branch distance of ≤ 3 to indicate a short distance or macro-



consistency (less than half of the total 7 branches), whereas a branch distance of  $> 3$  was regarded as indicative of inconsistency in the gene age.

## **The comparison of conflicting genes between GenTree and GageTracker**

To compare the accuracy of conflicting genes between GageTracker and GenTree, we used a similar method proposed by a previous study (14). This method involved the support evidence from three ortholog databases: OrthoDB, FlyBase, and Ensembl. First, we retrieved the orthologous gene clusters of *Drosophila* (level 7214 and species 7214 in December 2022) from OrthoDB by our self-written crawler script ([https://github.com/RiversDong/tools/blob/main/download\\_data\\_orthoDB.py](https://github.com/RiversDong/tools/blob/main/download_data_orthoDB.py)), FlyBase ([https://ftp.flybase.org/releases/FB2022\\_01/precomputed\\_files/orthologs/](https://ftp.flybase.org/releases/FB2022_01/precomputed_files/orthologs/)), and Ensembl biomaart (<https://dec2021-metazoa.ensembl.org/biomart/martview/>). Subsequently, we re-inferred the age of conflicting genes (between GageTracker and GenTree) using OrthoDB, FlyBase, and Ensembl, respectively with the same voting cutoff that was employed by GageTracker (voting=0.5). Afterward, we counted the number of consistency genes for differential gene classes and calculated the support ratio based on the branch distance to estimate the consistency between GageTracker and GenTree. The difference value here refers to the absolute age difference between genes. For instance, if GageTracker places a gene at branch 6, while GenTree estimates its age as branch 5, then the difference in age between the two is 1 ( $|\text{Age in GageTracker} - \text{Age in GenTree}|$ ). A smaller difference value indicates greater consistency between the annotation results.

## **Expression data of *D. melanogaster***

FlyAtlas 2.0 (<https://motif.mvls.gla.ac.uk/FlyAtlas2/>) was an online database that provided expression data for male, female, and larval stages of *D. melanogaster* (38). To investigate the expression patterns of genes with conflicting gene ages, we downloaded (in May 2023) gene expression data for 15 male tissues and 16 female tissues from FlyAtlas 2.0. Specifically, we focused on two groups of genes with conflicting age annotations, which we labeled as group A and group B. Group A comprised genes annotated as belonging to branch 0 by GageTracker but annotated as new genes (ranging from branch 1 to branch 6) by GenTree, while group B included genes annotated as old genes by GenTree but as new genes by GageTracker.

251

## 252 **RESULTS**

### 253 **Features and implementation of GageTracker**

254 GageTracker was an available software that employed a single command line execution to estimate  
 255 gene age. Figure 1A provided a visual representation of the software's architecture. The input data  
 256 required by GageTracker was only a GTF annotation of the focal species and the FASTA files for each  
 257 genome assembly. The primary objective of this algorithm was to determine whether a gene in the  
 258 focal species had orthologous counterparts in the outgroup species (middle panel of Figure 1B). The  
 259 algorithm aimed to identify a set of optimal symmetric alignments (micro-synteny) between the focal  
 260 and query genomes and merge these alignments into larger synteny block regions (macro-synteny).  
 261 The presence of a gene in the query species was indicated if it overlapped with a micro-synteny region  
 262 in both the focal and query genomes. Conversely, the absence of the gene in the outgroup species was  
 263 inferred if it failed to overlap with any micro-synteny region in the outgroup genome. To minimize  
 264 false positives in complex regions such as sequence gaps, centromeres, and telomeres, the gene age  
 265 estimation was limited to genes within macro-synteny blocks.

266 In summary, the GageTracker software consisted of five key processes:

267 1) Soft masking of repeat sequences in both the reference and query genomes using either the  
 268 WindowMasker or Tantan software (41,42), then converting the exon regions to uppercase if they  
 269 have been masked by WindowMasker or Tantan.

270 2) Identification of a set of local pairwise alignments between the focal and query genomes  
 271 through the seed-and-extend heuristic with the Last aligner (seeds were excluded from overlapping  
 272 masked repeat sequences, but the final alignments were allowed to extend into them).

273 3) Retrieval of Reciprocal Best Syntenic Net hits (RBH) using the UCSC Genome Browser  
 274 toolkit and conducting micro-synteny analysis between the two genomes.

275 4) Construction of genome alignment blocks, conducting macro-synteny analysis, and removal of  
 276 genes in genome alignment gaps with low confidence.

277 5) Combination of the genes' presence/absence status in all species and tracing of gene age  
 278 through a voting algorithm.

279

## Appropriate repeat-masking makes GageTracker more time-efficient

The eukaryotic genome was composed of a vast number of transposable elements (TEs) in various stages of evolution. TEs generated intricate repeat patterns, contributing up to 44% of the human genome (43), 50% of mouse genomes (44), and 64% of the maize genome (45,46). Masked repeats (such as interspersed repeats and simple sequences) were a common pre-alignment step to simplify genome alignment and reduce excessive output. However, excessive repeat-masking posed a risk to ortholog search which would also lead to false positives. If the true ortholog of a genome was masked while a paralog was not, the paralog would be aligned incorrectly, leading to an inaccurate result. To solve this, we integrated an optimized method for filtering repeated sequences. As GageTracker was primarily designed to handle various focal species, it would be beneficial to use a software that could predict repeat sequences from scratch rather than relying solely on the repeatMasker database. We used Windowmasker/Tantan to soft mark the repeat sequence from scratch, revised the wrongly masked sequence in the protein-coding region, and then used the seeds-extend algorithm to search genome alignments. This could minimize aligning dissimilar and unrelated sequences, as if two dissimilar genes had the similar flanking sequences, the overall alignment score was likely to exceed the threshold value and produce a false-positive.

To assess the improvement in running time for genome alignments using the new soft-masked strategy, we conducted tests on a larger simulated dataset with distant genomic divergence (“simMammals” dataset in Alignathon (37)). The query-reference pairs (simCow-simHuman, simDog-simHuman, simMouse-simHuman, simRat-simHuman) were aligned using Lastal and GageTracker strategy (Lastal: last+without repeat-masking; GageTracker: last+repeat masking, wrapped in Gagetacker -step1).

GageTracker masked approximately 8.73% of the simHuman genome and 5% of the other four simulated datasets (Table 1). The genome alignment results demonstrated that GageTracker performed outperformed Lastal in terms of speed, with a reduction in running time of 18-22% ( $p\text{-value} < 0.01$ ) by introducing repeat-masking strategy. Furthermore, GageTracker significantly reduced the output alignments (70-75%,  $p\text{-value} < 0.01$ ) compared to Lastal, thereby could result in substantial time savings in subsequent RBH steps.

Moreover, we also evaluated the time cost of GageTracker based on several pairs of real datasets, including *D. melanogaster* vs. *D. yakuba* and *H. sapiens* vs. *M. musculus*. As expected, the real

310 datasets had higher levels of repeat sequences. GageTracker masked approximately 22.3% of the  
 311 repeat sequences in the *D. melanogaster*, leading to a significant reduction of 43.3% in calculations  
 312 and a 46.7% saving in running time than Lastal. Similarly, in the *H. sapiens* vs. *M. musculus* genomes,  
 313 GageTracker masked 36.2% of repeat sequences in the *H. sapiens* genome, reducing redundant  
 314 calculations by 60.3% and saving 68.7% of the runtime than Lastal. Together, these results indicated  
 315 that proper repeat-masking in GageTracker could significantly save running time and reduce the  
 316 output in genome alignment.

Table 1. Genome alignment between Lastal, GageTracker and LastZ.

Genome	Genome size	True repeat	Repeat mask (GageTracker)	MAF file			Runtime (min)		
				Lastal	GageTracker	LastZ	Lastal	GageTracker	LastZ
simHuman	191 Mb	11.09%	8.73%	-	-	-			
simCow	193 Mb	9.88%	4.85%	12.27 Gb	3.61 Gb	2.45 Gb	425	323	1958
simDog	192 Mb	10.18%	5.08%	12.52 Gb	3.71 Gb	2.24 Gb	356	267	2324
simMouse	199 Mb	12.13%	4.84%	8.42 Gb	2.13 Gb	0.64 Gb	363	281	445
simRat	199 Mb	11.47%	5.04%	8.26 Gb	2.05 Gb	0.68 Gb	436	356	541
<i>D. melanogaster</i>	144 Mb	-	22.32%	-	-	-	-	-	-
<i>D. yakuba</i>	166 Mb	-	7.80%	10.06 Gb	5.71 Gb	-	197	105	-
<i>H. sapiens</i>	3.09 Gb	-	36.21%	-	-	-	-	-	-
<i>M. musculus</i>	2.72 Gb	-	7.51%	1068.83 Gb	424.49 Gb	-	5412	1692	-

Note. These query-reference pairs (simCow-simHuman, simDog-simHuman, simMouse-simHuman, simRat-simHuman, *D. melanogaster* - *D. yakuba* and *H. sapiens* - *M. musculus* a) were aligned with Lastal, GageTracker, and LastZ aligner, respectively. The bold species represent the reference genome. Lastal: last + without repeat-masking; GageTracker: last + repeat-masking; LastZ: LastZ + truth repeats. When performing LastZ without repeat-masking, it failed to finish the alignment in 7 days and is excluded from the table. All runs are controlled in 6 threads on the servers with an E7-4850-CPU (2.5GHZ) array and 1 TB of memory.

## 326 **GageTracker shows high precision and recall in ortholog search**

327 To determine whether repeat-masking would impair the quality of the GageTracker, we  
 328 compared it with LastZ. LastZ was a widely used and powerful genome alignment software  
 329 that was relied upon by other software such as UCSC, Cactus, and SegAlign (47-49). The  
 330 LastZ with suggested parameters (O=400 E=30 L=3000 H=2200 T=1) in UCSC could  
 331 achieve high sensitivity but also result in much more time-consuming. To address this, we  
 332 used PLastZ (<https://github.com/AntoineHo/PLastZ>), a parallel LastZ script, to improve  
 333 alignment speed. Then, we used LastZ as the benchmark and compared the accuracy from  
 334 these four strategies: Lastal (last + without repeat-masking), GageTracker (last + repeat-  
 335 masking), LASTZ (LastZ + without repeat-masking), and LastZ (LastZ + repeat-masking).  
 336 Generally, assessing the quality of genome aligners was challenging because the true  
 337 homologous sequences were unknown. The Alignathon simulated datasets provided a solution,  
 338 as they included a known truth data set for evaluating multiple genome aligners (50-53). To  
 339 search for one-to-one orthologous alignments, we used the no paralog simulated Mammals  
 340 data as the truth and the simHuman genome as a reference. We applied chains-and-nets and  
 341 Syntenic net algorithms to find reciprocal best alignments from test data generated by Lastal,  
 342 GageTracker, LASTZ, and LastZ strategies and employed mafTransitive and Closure  
 343 mafComparator to calculate precision, recall, and *F*-score values. All runs were controlled in  
 344 6 threads on the servers with E7-4850-CPU (2.5GHZ) array respectively.

345 Without using repeat-masking strategy, LASTZ (LastZ + without repeat-masking) could  
 346 not complete genome alignments within 7 days in the simulated data, far exceeding the  
 347 expected time and making them impractical. But Lastal (last + without repeat-masking)  
 348 finished genome alignments in 356~436 minutes (Table 1) and reciprocal best alignments in  
 349 430~513 minutes (Figure 1C), indicating that Lastal ran much faster than LASTZ. When  
 350 using repeat-masking process, LastZ employed true repeat sequences and masked a high  
 351 proportion of repeats (~10%), which should theoretically yield the fastest and most accurate  
 352 data (Table 1). Nevertheless, the result showed that GageTracker had a faster computation  
 353 speed and reduced runtime by 1.4-7 times compared to LastZ both in genome alignment and  
 354 reciprocal best alignments process (Table1 and Figure 1C). Based on our analysis,  
 355 GageTracker was found to be the fastest in terms of computational speed, followed by Lastal,

356 LastZ, and LASTZ, in descending order.

357 In addition, GageTracker and LastZ achieved comparable precision, recall, and *F-score*  
 358 values across four mammal alignment groups (simHuman vs. simRat, simHuman vs.  
 359 simMouse, simHuman vs. simDog, simHuman vs. simCow, Figure 1D). It also slightly  
 360 improved *F-scores* and runtime compared to Lastal. These findings suggested that, in terms of  
 361 calculation accuracy, GageTracker performed similarly to LastZ but slightly higher than  
 362 Lastal. Thus, GageTracker was suitable for searching orthologous fragments, offering high  
 363 sensitivity, accuracy, and efficient processing.

364

### 365 **Gene age exhibits a high level of consistency in the result of GageTracker and GenTree** 366 **in the real fly data set**

367 We selected *Drosophila melanogaster*, which was an important model organism in genetic  
 368 research with closely related species of improved sequencing quality (Middle panel of Figure  
 369 1B), to compare GageTracker and GenTree consistency in three levels (please refer to our  
 370 method section: Age comparison in three different levels). GenTree was a reliable gene age  
 371 database that integrated synteny-based, protein-family-based, and small-scale manual curation  
 372 methods with functional genomic data to infer gene age (4). GenTree was built as a central  
 373 portal to catalog lineage-specific genes in several model species. By using GageTracker, it  
 374 took about 22 hours to estimate 23,720 genes' age (including ncRNA and protein genes) in 12  
 375 *Drosophila* sequencing genomes, of which approximately 13,965 were protein genes  
 376 (Supplementary Table S2). Note that the two methods used different annotations, which  
 377 would lead to a difference in gene number. We selected the overlapping protein genes in both  
 378 annotations as our comparison objects (12,888, Figure 1E) to compare GageTracker and  
 379 GenTree. On the global level, GageTracker exhibited high consistency with GenTree, with  
 380 ~94.4% of the genes (12,171/12,888) sharing the same gene age (Figure 1F and  
 381 Supplementary Table S3). Only 5.6% of the genes showed conflicts in exact gene age  
 382 between the two methods (717, 2.4%+3.2%=5.6%, Figure 1F, and Supplementary Table S3).  
 383 Additionally, when comparing the old and young gene sets (where branch 0 was considered  
 384 old and branches 1-6 were considered young), GageTracker showed approximately 98%  
 385 consistency with GenTree in the old catalog and approximately 80% consistency in the young

catalog (Figure 1G). At a detailed level, it was found that genes sharing the same age annotation accounted for more than 70% of certain branches (over 98% in branch 0 and almost 80% for genes in branches 5 and 6), with the exception of branches 2 and 3 (Figure 1E). However, for branches 2 and 3, the consistency still approached 50%. The consistency of the gene age presented in this work was quite high if we made a parallel comparison. For example, Shao et al. ever compared their results (GenTree) with the results from the other two previous results, and they showed that there were approximately 30% and 40% conflicts between GenTree and the other two related works, respectively. However, the comparison between GageTracker and GenTree on the global level showed a high level of consistency.

Additionally, it was important to note that the conflicting gene age in some branches could be caused by differences in genome assembled versions and gene-age dating methods between GageTracker and GenTree. To evaluate the deviation between the two methods, we used the branch distance metric. The results show that more than 97.94% of genes have a branch distance of no more than 2, and more than 99% of genes have a branch distance of no more than 3 (half of the branches are 3), indicating that the age annotated by the two methods was very consistent. Together, these results suggested that GageTracker and GenTree achieved a high level of consistency at global level and macroscopic level.

403

#### 404 **GageTracker avoids overestimating the gene age on conflicting genes**

The inability to detect alignment hits in sequencing gaps was addressed by introducing macro-synteny in our study, but GenTree did not consider this issue. To further clarify the accuracy of GageTracker, we checked the ortholog genes of these 717 conflicting genes from three different ortholog databases: OrthoDB, Ensembl, and FlyBase. The corresponding orthologous genes of these 717 genes in 12 *Drosophila* genomes were screened out and the genes' age was also estimated by the same voting cutoff that was used in GageTracker.

In OrthoDB, 667 conflicting genes had corresponding orthologs in 12 *Drosophila* species (Supplementary Table S4). Approximately 36.28% (242/667) of conflicting genes had the same gene ages between GageTracker vs. OrthoDB, whereas 32.53% (217/667) were consistent between GenTree vs. OrthoDB, exhibiting that the consistency between GageTracker vs. OrthoDB was slightly higher than that between GenTree vs. OrthoDB. This



416 result also suggested that a considerable proportion of the conflicting genes shared the same  
417 gene age annotation between either GageTracker and OrthoDB or GenTree and OrthoDB. In  
418 addition, we evaluated the conflicting genes by comparing the gene numbers of each branch  
419 distance between GageTracker vs. OrthoDB, and GenTree vs. OrthoDB. A higher proportion  
420 of genes in GageTracker (80.06%) had a shorter branch distance with OrthoDB (distance  $\leq 2$   
421 branches,  $p\text{-value}=0.0137$ , Figure 2A) compared to GenTree (67.31%), which suggested that  
422 GageTracker was more consistent with OrthoDB. On the other hand, GenTree had a higher  
423 proportion of genes with a larger branch distance (distance  $\geq 3$  branches,  $p\text{-value} = 0.0075$ ),  
424 indicating a greater propensity for overestimating gene age compared to GageTracker.

425 In Ensembl, we observed 717 conflicting genes that had corresponding annotations in  
426 our downloaded orthologs (Supplementary Table S5). We compared the gene numbers of each  
427 branch distance between GageTracker vs. Ensembl, and GenTree vs. Ensembl group.  
428 GageTracker vs. Ensembl group still showed a higher proportion gene with a shorter branch  
429 distance (0-1), while GenTree vs. Ensembl group had a higher proportion gene with a larger  
430 branch distance (Figure 2B, branch distance of 2-6,  $p\text{-value} = 0.0021$ ). This result also  
431 indicated that, compared to Gentree, GageTracker showed better consistency with Ensembl.

432 In FlyBase, 665 conflicting genes traced corresponding orthologs in 12 *Drosophila*  
433 species (Supplementary Table S6). GenTree vs. FlyBase group also exhibited a higher  
434 proportion of genes with larger branch distance (Figure 2C, branch distance of 3-6,  $p\text{-value}$   
435  $=0.0266$ ) than the GageTracker vs. FlyBase group. To summarize, for the conflicting age,  
436 GageTracker avoided overestimation of genes age and obtained a more consistent with three  
437 other homologous sources (OrthoDB, FlyBase, and Ensembl), compared to GenTree. Since  
438 OrthoDB, Ensembl, and Flybase could not be taken as ‘golden age’ datasets, these analyses at  
439 least suggested that GageTracker's accuracy could be comparable to that of GenTree.

440 To investigate why GageTracker and GenTree showed age differences compared to three  
441 orthologous databases, we examined the copy numbers of these conflicting genes in their  
442 corresponding gene family (Supplementary Table S7) based on the paralogs recorded in  
443 FlyBase (dmel\_paralogs\_fb\_2022\_03.tsv.gz). It was found that approximately 62.12%  
444 (443/717) of conflicting genes were multiple copy genes and ~45% had more than two copies  
445 (Figure 2D). For conflicting genes with a larger branch distance ( $>3$ ), around ~72% of them

were found to be multiple copy genes, indicating that the conflicting genes with larger evolutionary branch bias were likely due to more paralogs (Figure 2D). The high sequence similarity of paralogous sequences increased the complexity of the gene family by adding more copies and/or changing their collinearity, which could result in more difficulty to infer the gene age accurately. Considering the result from the analysis of three ortholog databases, GageTracker demonstrates a higher level of accuracy in estimating gene age.

### **The expression of conflicting genes indicates that young genes in GageTracker but old in GenTree have a preference expression in testis**

Previous research has shown that old genes tended to be expressed more broadly, whereas new genes were more tissue-specific and enriched in testes (3,39,54). Therefore, we conducted an expression pattern analysis on a subset of genes that showed conflicting gene age. First, we divided the genes into two groups (designated A and B) based on their age annotations in GageTracker and GenTree. The genes annotated as belonging to branch 0 in GageTracker and branches 1 to 6 in GenTree formed Group A, while the genes annotated as belonging to branches 1 to 6 in GageTracker and branch 0 in GenTree formed Group B. Using gene expression data from FlyAtlas (38), we compared the expression patterns of these two conflicting age groups. Our analysis revealed that, in male, group A genes demonstrated higher average expression levels than group B genes in 11 out of 15 tissues (Figure 3, Supplementary Table S8, over 73% of tissues). Additionally, group B genes exhibited a significant bias towards expression in the testis (Figure 3, Wilcoxon rank-sum test,  $p$ -value =  $5.48 \times 10^{-11}$ ). In females, group A genes were also expressed at higher mean levels than group B genes in 10 different tissues (Figure 4, Supplementary Table S9). Overall, our findings suggested that genes classified as old by GageTracker exhibited higher average expression levels and a wider range of expression, while genes categorized as young by GageTracker but labeled as old by GenTree displayed a significant tendency to be expressed in the testis. To summarize, for genes with conflicting age annotations, those identified as old genes by GageTracker showed expression patterns typical of old genes, while those identified as young genes displayed male-biased expression patterns characteristic of young genes.

## 476 DISCUSSION

477 GageTracker was a software tool that has two operation models: wrapped running mode and  
 478 step-by-step running mode. The wrapped running mode integrated all necessary steps for gene  
 479 age annotation, such as genome alignment, RBH calculation, and age assignment, into a  
 480 single process specified in a configuration file. This mode was convenient for users as it  
 481 separated command parameters and datasets. However, genome alignment and RBH  
 482 calculation were the most time-consuming stages, and users may want to add new branches  
 483 and outgroup species to existing genome alignment results. To address this, the step-by-step  
 484 running mode was introduced, allowing users to add new genome alignments without  
 485 affecting previous results, thus could save running time. The step-by-step mode also enabled  
 486 the addition of new branches and focal species. In conclusion, GageTracker facilitated the  
 487 annotation of new genes through its two operation modes (the wrapped mode integrates all  
 488 necessary steps into a single process, while the step-by-step mode allows for the addition of  
 489 new genome alignments without affecting previous results). GageTracker also integrated an  
 490 RBH detection pipeline that was used by UCSC Genome Browser, which could help  
 491 GageTracker find reciprocal best regions in a more accurate and more fast manner than other  
 492 genome-based aligners (34).

493 Several reasons could affect the estimation of new genes, for example, the quality of  
 494 genome assembly, genetic relationship between focal and query species, and sequence  
 495 complexity, among which including high-quality genomes of closely related species help to  
 496 accurately estimate the number of new genes. However, due to the lack of high-quality  
 497 assembled genomes for some closely related species, new gene numbers may have been  
 498 overestimated. Specifically, a gene may have been annotated as new not because it is absent  
 499 in focal species, but rather due to a sequencing gap. As a result, we suggested that users could  
 500 select a more complete reference genome to mitigate false positives resulting from genome  
 501 sequencing gaps. By doing so, a more precise assessment of the number of novel genes could  
 502 be achieved. Additionally, we recommend selecting 2-3 species from each branch to prevent  
 503 false positives caused by gene loss. Furthermore, we have incorporated a criterion known as  
 504 macro-synteny into GageTracker to identify sequencing gaps. Any genes located in such  
 505 regions will be excluded from the annotation procedure since new genes in such regions may

506 be false positives.

507 It should be noted that GageTracker, in its current version, was unable to differentiate  
 508 between protein-coding genes and orthologous noncoding sequences. This distinction was a  
 509 complex process that required comprehensive genomic and expression data to determine  
 510 whether highly similar orthologous sequences in outgroup species were noncoding sequences.  
 511 GenTree, on the other hand, utilized both evolutionary and proteomic evidence to assess the  
 512 protein-coding potential of a sequence and overcome this limitation. Due to its inability to  
 513 distinguish between protein-coding genes and orthologous noncoding sequences,  
 514 GageTracker may overestimate the age of some young genes, leading to a slight  
 515 overestimation of their true age (the estimated age is older than the actual age). However, for  
 516 young genes, GageTracker was still able to provide reliable new gene candidates and their  
 517 estimated origin age (i.e., the age when homologous sequences first emerged). This allowed  
 518 researchers to trace back from the estimated ancient origin age and identify evidence of the  
 519 gradual evolution of orthologous noncoding sequences into *de novo* genes. We have recently  
 520 developed a random forest-based model that can detect coding capacity (55), which provided  
 521 a means for estimating protein-coding potential. In the upcoming version of GageTracker, we  
 522 plan to integrate the algorithm with transcriptome data to enable more accurate identification  
 523 of gene structures. This will allow us to perform a more detailed analysis of the origin and  
 524 evolutionary process of young genes, and generate structural diagrams for better visualization.  
 525 This integration will improve the accuracy of our predictions and provide a more  
 526 comprehensive understanding of the evolutionary history of young genes.

527

528

## 529 SUPPLEMENTARY DATA

530 Table S1. The genomes used in this study.

531 Table S2. Protein-coding genes' age of *Drosophila melanogaster* dated by GageTracker.

532 Seven genes are located outside the block region, with low reliability (NCON)

533 Table S3. The estimation of gene age between GageTracker and GenTree.

534 Table S4. The gene age inferred by orthoDB for the conflicting genes between GageTracker

535 and GenTree.

536 Table S5. The conflicting gene age inferred by Ensembl orthologs database between  
537 GageTracker and GenTree.

538 Table S6. The gene age inferred by the Flybase orthologs database for the conflicting genes  
539 between GageTracker and GenTree.

540 Table S7. Paralogous of conflicting genes from the Flybase homologs database.

541 Table S8. Expression pattern of conflicting genes in male *Drosophila melanogaster*.

542 Table S9. Expression pattern of conflicting genes in female *Drosophila melanogaster*.

543

## 544 **AUTHOR CONTRIBUTIONS**

545 SH, HZ, CF, and CD conceived and designed the study. CF and CD collected the data and  
546 performed the analysis. CW, FX, HZ, LY, XG, SL, WF and TL provided technical support. CF  
547 and CD wrote the manuscript. SH and XG revised the paper. All authors have read and  
548 accepted the final version of the manuscript.

549

550

## 551 **Acknowledgments**

552 We would like to thank Dr. Suxiang Lu, Ning Sun, and Han Xu for testing the GageTracker  
553 and providing some valuable suggestions.

554

555

## 556 **Funding**

557 This project is supported by the “Strategic Priority Research Program” of the Chinese  
558 Academy of Science (Grant No. XDB31000000), the National Natural Science Foundation of  
559 China (Grant, No. 32170438).

560

561

## 562 **CONFLICT OF INTEREST**

563 The authors declare that they have no known competing financial interests or personal  
564 relationships that could have appeared to influence the work reported in this paper.

565

566

567

568

569

## REFERENCES

1. Chen, S., Krinsky, B.H. and Long, M. (2013) New genes as drivers of phenotypic evolution. *Nat Rev Genet*, **14**, 645-660.
2. Sayers, E.W., Bolton, E.E., Brister, J.R., Canese, K., Chan, J., Comeau, D.C., Farrell, C.M., Feldgarden, M., Fine, A.M., Funk, K. *et al.* (2023) Database resources of the National Center for Biotechnology Information in 2023. *Nucleic Acids Res*, **51**, D29-D38.
3. Zhang, W., Landback, P., Gschwend, A.R., Shen, B. and Long, M. (2015) New genes drive the evolution of gene interaction networks in the human and mouse genomes. *Genome Biol*, **16**, 202.
4. Shao, Y., Chen, C., Shen, H., He, B.Z., Yu, D., Jiang, S., Zhao, S., Gao, Z., Zhu, Z., Chen, X. *et al.* (2019) GenTree, an integrated resource for analyzing the evolution and function of primate-specific coding genes. *Genome Res*, **29**, 682-696.
5. Zhang, Y.E., Landback, P., Vibranovski, M.D. and Long, M. (2011) Accelerated recruitment of new brain development genes into the human genome. *PLoS Biol*, **9**, e1001179.
6. Long, M.Y., VanKuren, N.W., Chen, S.D. and Vibranovski, M.D. (2013) New Gene Evolution: Little Did We Know. *Annu Rev Genet*, **47**, 307-333.
7. Levy, A. (2019) Genes from the Junkyard. *Nature*, **574**, 314-316.
8. Wei, W., Jin, Y.T., Du, M.Z., Wang, J., Rao, N. and Guo, F.B. (2016) Genomic Complexity Places Less Restrictions on the Evolution of Young Coexpression Networks than Protein-Protein Interactions. *Genome Biol Evol*, **8**, 2624-2631.
9. Fang, C., Gan, X., Zhang, C. and He, S. (2021) The new chimeric chiron genes evolved essential roles in zebrafish embryonic development by regulating NAD(+) levels. *Sci China Life Sci*, **64**, 1929-1948.
10. VanKuren, N.W. and Long, M. (2018) Gene duplicates resolving sexual conflict rapidly evolved essential gametogenesis functions. *Nat Ecol Evol*, **2**, 705-712.
11. Long, M. and Langley, C.H. (1993) Natural selection and the origin of jingwei, a chimeric processed functional gene in *Drosophila*. *Science*, **260**, 91-95.
12. Dai, H., Chen, Y., Chen, S., Mao, Q., Kennedy, D., Landback, P., Eyre-Walker, A., Du, W. and Long, M. (2008) The evolution of courtship behaviors through the origination of a new gene in *Drosophila*. *Proc Natl Acad Sci U S A*, **105**, 7478-7483.
13. Chen, S., Zhang, Y.E. and Long, M. (2010) New genes in *Drosophila* quickly become essential. *Science*, **330**, 1682-1685.
14. Xia, S.Q., VanKuren, N.W., Chen, C.Y., Zhang, L., Kemkemer, C., Shao, Y., Jia, H.X., Lee, U., Advani, A.S., Gschwend, A. *et al.* (2021) Genomic analyses of new genes and their phenotypic effects reveal rapid evolution of essential functions in *Drosophila* development. *Plos Genet*, **17**.
15. Xia, S., Wang, Z., Zhang, H., Hu, K., Zhang, Z., Qin, M., Dun, X., Yi, B., Wen, J., Ma, C. *et al.* (2016) Altered Transcription and Neofunctionalization of Duplicated Genes Rescue the Harmful Effects of a Chimeric Gene in *Brassica napus*. *Plant Cell*, **28**, 2060-2078.
16. Huang, Y., Chen, J., Dong, C., Sosa, D., Xia, S., Ouyang, Y., Fan, C., Li, D., Mortola, E., Long, M. *et al.* (2022) Species-specific partial gene duplication in *Arabidopsis thaliana* evolved novel phenotypic effects on morphological traits under strong positive selection. *Plant Cell*, **34**, 802-817.
17. Zhou, Y., Zhang, C., Zhang, L., Ye, Q., Liu, N., Wang, M., Long, G., Fan, W., Long, M. and

- Wing, R.A. (2022) Gene fusion as an important mechanism to generate new genes in the genus *Oryza*. *Genome Biol*, **23**, 130.
18. Zhang, L., Ren, Y., Yang, T., Li, G., Chen, J., Gschwend, A.R., Yu, Y., Hou, G., Zi, J., Zhou, R. *et al.* (2019) Rapid evolution of protein diversity by de novo origination in *Oryza*. *Nat Ecol Evol*, **3**, 679-690.
19. Jin, G., Ma, P.F., Wu, X., Gu, L., Long, M., Zhang, C. and Li, D.Z. (2021) New Genes Interacted With Recent Whole-Genome Duplicates in the Fast Stem Growth of Bamboos. *Mol Biol Evol*, **38**, 5752-5768.
20. Wei, W., Zhang, T., Lin, D., Yang, Z.J. and Guo, F.B. (2013) Transcriptional abundance is not the single force driving the evolution of bacterial proteins. *BMC Evol Biol*, **13**, 162.
21. Wei, W., Ning, L.W., Ye, Y.N. and Guo, F.B. (2013) Geptop: a gene essentiality prediction tool for sequenced bacterial genomes based on orthology and phylogeny. *PLoS One*, **8**, e72343.
22. Ezkurdia, I., Juan, D., Rodriguez, J.M., Frankish, A., Diekhans, M., Harrow, J., Vazquez, J., Valencia, A. and Tress, M.L. (2014) Multiple evidence strands suggest that there may be as few as 19,000 human protein-coding genes. *Hum Mol Genet*, **23**, 5866-5878.
23. Toll-Riera, M., Bosch, N., Bellora, N., Castelo, R., Armengol, L., Estivill, X. and Alba, M.M. (2009) Origin of primate orphan genes: a comparative genomics approach. *Mol Biol Evol*, **26**, 603-612.
24. Tong, Y.B., Shi, M.W., Qian, S.H., Chen, Y.J., Luo, Z.H., Tu, Y.X., Xiong, Y.L., Geng, Y.J., Chen, C. and Chen, Z.X. (2021) GenOrigin: A comprehensive protein-coding gene origination database on the evolutionary timescale of life. *J Genet Genomics*, **48**, 1122-1129.
25. Tay, S.K., Blythe, J. and Lipovich, L. (2009) Global discovery of primate-specific genes in the human genome. *Proc Natl Acad Sci U S A*, **106**, 12019-12024.
26. Zhang, Y.E., Vibrantovski, M.D., Landback, P., Marais, G.A. and Long, M. (2010) Chromosomal redistribution of male-biased genes in mammalian evolution with two bursts of gene gain on the X chromosome. *PLoS Biol*, **8**, e1000494.
27. Stein, J.C., Yu, Y., Copetti, D., Zwickl, D.J., Zhang, L., Zhang, C., Chougule, K., Gao, D., Iwata, A., Goicoechea, J.L. *et al.* (2018) Genomes of 13 domesticated and wild rice relatives highlight genetic conservation, turnover and innovation across the genus *Oryza*. *Nat Genet*, **50**, 285-296.
28. Kabza, M., Ciomborowska, J. and Makalowska, I. (2014) RetrogeneDB--a database of animal retrogenes. *Mol Biol Evol*, **31**, 1646-1648.
29. Kriventseva, E.V., Kuznetsov, D., Tegenfeldt, F., Manni, M., Dias, R., Simao, F.A. and Zdobnov, E.M. (2019) OrthoDB v10: sampling the diversity of animal, plant, fungal, protist, bacterial and viral genomes for evolutionary and functional annotations of orthologs. *Nucleic Acids Res*, **47**, D807-D811.
30. Kuznetsov, D., Tegenfeldt, F., Manni, M., Seppey, M., Berkeley, M., Kriventseva, E.V. and Zdobnov, E.M. (2023) OrthoDB v11: annotation of orthologs in the widest sampling of organismal diversity. *Nucleic Acids Res*, **51**, D445-D451.
31. Fan, G., Song, Y., Yang, L., Huang, X., Zhang, S., Zhang, M., Yang, X., Chang, Y., Zhang, H., Li, Y. *et al.* (2020) Initial data release and announcement of the 10,000 Fish Genomes Project (Fish10K). *Gigascience*, **9**, gaa080.
32. Zhang, G., Rahbek, C., Graves, G.R., Lei, F., Jarvis, E.D. and Gilbert, M.T. (2015) Genomics: Bird sequencing project takes off. *Nature*, **522**, 34.



33. Miao, W., Song, L., Ba, S., Zhang, L., Guan, G., Zhang, Z. and Ning, K. (2020) Protist 10,000 Genomes Project. *Innovation (Camb)*, **1**, 100058.
34. Frith, M.C. and Kawaguchi, R. (2015) Split-alignment of genomes finds orthologies more accurately. *Genome Biol*, **16**, 106.
35. Kolmogorov, M., Armstrong, J., Raney, B.J., Streeter, I., Dunn, M., Yang, F., Odom, D., Flicek, P., Keane, T.M., Thybert, D. *et al.* (2018) Chromosome assembly of large and complex genomes using multiple references. *Genome Res*, **28**, 1720-1732.
36. Quinlan, A.R. and Hall, I.M. (2010) BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*, **26**, 841-842.
37. Earl, D., Nguyen, N., Hickey, G., Harris, R.S., Fitzgerald, S., Beal, K., Seledtsov, I., Molodtsov, V., Raney, B.J., Clawson, H. *et al.* (2014) Alignathon: a competitive assessment of whole-genome alignment methods. *Genome Res*, **24**, 2077-2089.
38. Leader, D.P., Krause, S.A., Pandit, A., Davies, S.A. and Dow, J.A.T. (2018) FlyAtlas 2: a new version of the *Drosophila melanogaster* expression atlas with RNA-Seq, miRNA-Seq and sex-specific data. *Nucleic Acids Res*, **46**, D809-D815.
39. Zhang, Y.E., Vibranovski, M.D., Krinsky, B.H. and Long, M. (2010) Age-dependent chromosomal distribution of male-biased genes in *Drosophila*. *Genome Res*, **20**, 1526-1533.
40. Fang, C., Zou, C., Fu, Y., Li, J., Li, Y., Ma, Y., Zhao, S. and Li, C. (2018) DNA methylation changes and evolution of RNA-based duplication in *Sus scrofa*: based on a two-step strategy. *Epigenomics*, **10**, 199-218.
41. Morgulis, A., Gertz, E.M., Schaffer, A.A. and Agarwala, R. (2006) WindowMasker: window-based masker for sequenced genomes. *Bioinformatics*, **22**, 134-141.
42. Frith, M.C. (2011) A new repeat-masking method enables specific detection of homologous sequences. *Nucleic Acids Res*, **39**, e23.
43. Mills, R.E., Bennett, E.A., Iskow, R.C. and Devine, S.E. (2007) Which transposable elements are active in the human genome? *Trends Genet*, **23**, 183-191.
44. Bruno, M., Mahgoub, M. and Macfarlan, T.S. (2019) The Arms Race Between KRAB-Zinc Finger Proteins and Endogenous Retroelements and Its Impact on Mammals. *Annu Rev Genet*, **53**, 393-416.
45. SanMiguel, P., Tikhonov, A., Jin, Y.K., Motchoulskaia, N., Zakharov, D., Melake-Berhan, A., Springer, P.S., Edwards, K.J., Lee, M., Avramova, Z. *et al.* (1996) Nested retrotransposons in the intergenic regions of the maize genome. *Science*, **274**, 765-768.
46. Jiao, Y., Peluso, P., Shi, J., Liang, T., Stitzer, M.C., Wang, B., Campbell, M.S., Stein, J.C., Wei, X., Chin, C.S. *et al.* (2017) Improved maize reference genome with single-molecule technologies. *Nature*, **546**, 524-527.
47. Paten, B., Earl, D., Nguyen, N., Diekhans, M., Zerbino, D. and Haussler, D. (2011) Cactus: Algorithms for genome multiple sequence alignment. *Genome Res*, **21**, 1512-1528.
48. Goenka, S.D., Turakhia, Y., Paten, B. and Horowitz, M. (2020) SegAlign: A Scalable GPU-Based Whole Genome Aligner. *Proceedings of Sc20: The International Conference for High Performance Computing, Networking, Storage and Analysis (Sc20)*.
49. Lee, B.T., Barber, G.P., Benet-Pages, A., Casper, J., Clawson, H., Diekhans, M., Fischer, C., Gonzalez, J.N., Hinrichs, A.S., Lee, C.M. *et al.* (2022) The UCSC Genome Browser database: 2022 update. *Nucleic Acids Res*, **50**, D1115-D1122.
50. Kille, B., Balaji, A., Sedlazeck, F.J., Nute, M. and Treangen, T.J. (2022) Multiple genome

- alignment in the telomere-to-telomere assembly era. *Genome Biol*, **23**, 182.
51. Minkin, I. and Medvedev, P. (2020) Scalable multiple whole-genome alignment and locally collinear block construction with SibeliaZ. *Nat Commun*, **11**, 6327.
52. Armstrong, J., Hickey, G., Diekhans, M., Fiddes, I.T., Novak, A.M., Deran, A., Fang, Q., Xie, D., Feng, S., Stiller, J. *et al.* (2020) Progressive Cactus is a multiple-genome aligner for the thousand-genome era. *Nature*, **587**, 246-251.
53. Armstrong, J., Fiddes, I.T., Diekhans, M. and Paten, B. (2019) Whole-Genome Alignment and Comparative Annotation. *Annu Rev Anim Biosci*, **7**, 41-64.
54. 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.
55. Yu, J., Jiang, W., Zhu, S.B., Liao, Z., Dou, X., Liu, J., Guo, F.B. and Dong, C. (2023) Prediction of protein-coding small ORFs in multi-species using integrated sequence-derived features and the random forest model. *Methods*, **210**, 10-19.

## TABLE AND FIGURES LEGENDS

Table 1. Genome alignment between Lastal, GageTracker and LastZ.

These query-reference pairs (simCow-simHuman, simDog-simHuman, simMouse-simHuman, simRat-simHuman, *D. melanogaster* - *D. yakuba* and *H. sapiens* - *M. musculus* a) were aligned with Lastal, GageTracker, and LastZ aligner. The bold species represent the reference genome. Lastal: last + without repeat-masking; GageTracker: last + repeat-masking; LastZ: LastZ + truth repeats. When performing LastZ without repeat-masking, it failed to finish the alignment in 7 days and is excluded from the table. All runs are controlled in 6 threads on the servers with an E7-4850-CPU (2.5GHZ) array and 1 TB of memory.

Figure 1. The pipeline and performance of GageTracker.

(A) The overview of inputting data (orange), UCSC Genome Browser toolkit (green), post-processing (purple), and output (red). (B) A diagram about alignment block (macro-synteny) and micro-synteny. The rectangles with the same colors mean the homologous genes (left panel in B). The left panel of B also illustrates the Voting strategy for assessing the support ratio of a gene's origin on a branch. The orange rectangle represents a gene that has an ortholog on the corresponding branch, while the white rectangle means absence. The “br” represents the short name of a branch. The number in the orange rectangle represents the support ratio of the gene originating in that branch. For example, Gene 1 in the dotted box has corresponding ortholog on branches 6, 3, 2, and 0. The support ratio of Gene 1's origin on branch 3 is 50% (2/4, present on 3 out of 4 branches), on branch 2 is 60% (3/5, present on 3 out of 5 branches), and on branch 0, its support rate reaches 57% (4/7). Then, this gene age was placed on its oldest branch 0 (red arrow). Gene 3 independently emerges on branch 6 as its support rate of originating on branch 0 is only 28.5%. The detailed calculation strategy could be seen in the Method. The middle panel of B is a phylogeny tree of *Drosophila* species and the blue lines indicated the evolution direction leading to *D. melanogaster*. The topology of the tree and the divergent time of each branch point were retrieved from the previous study (39). The right panel of B is the gene number comparison between GenTree and GageTracker. (C) the total running time of three aligners for producing reciprocal best alignments: Lastal

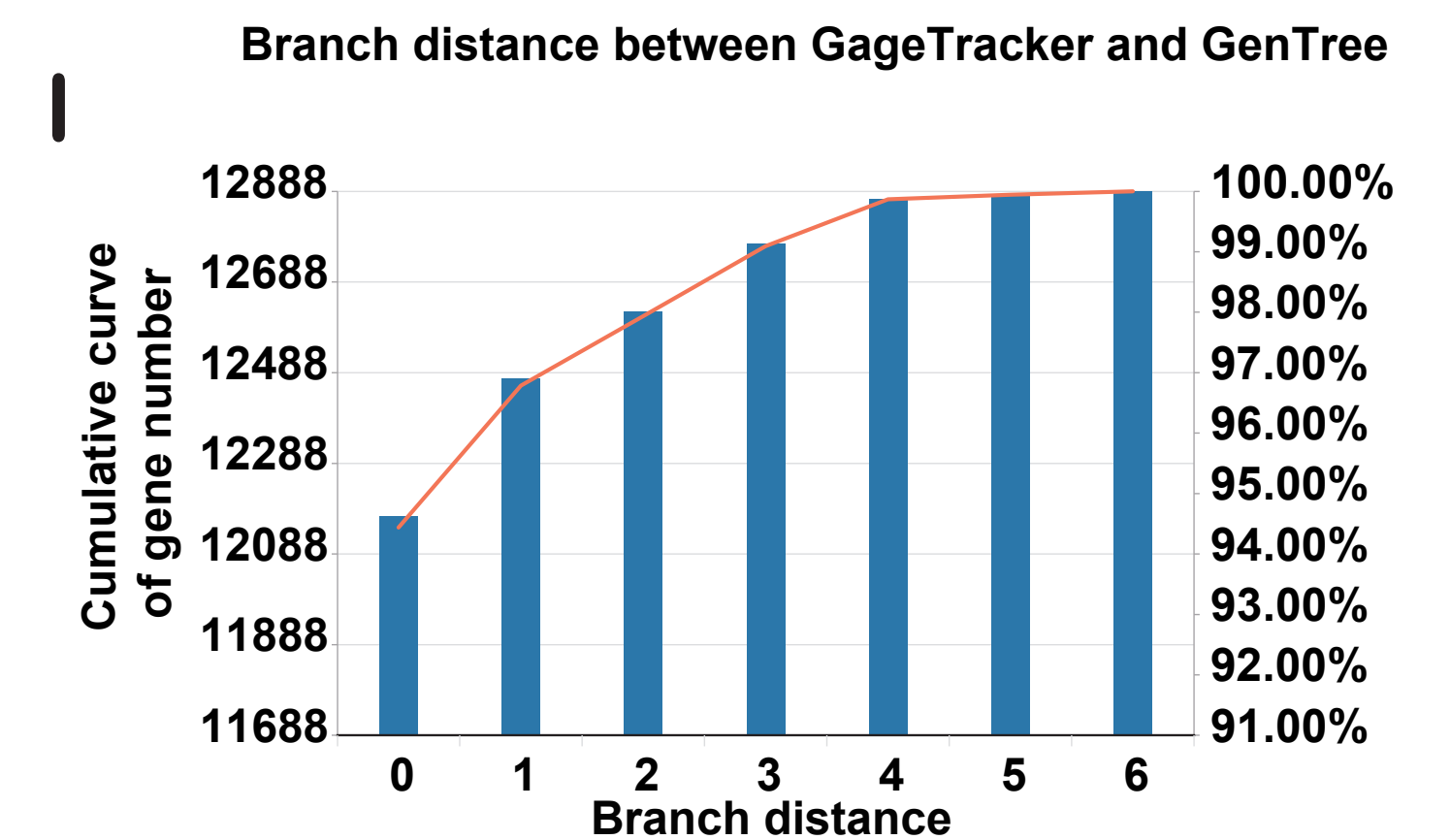
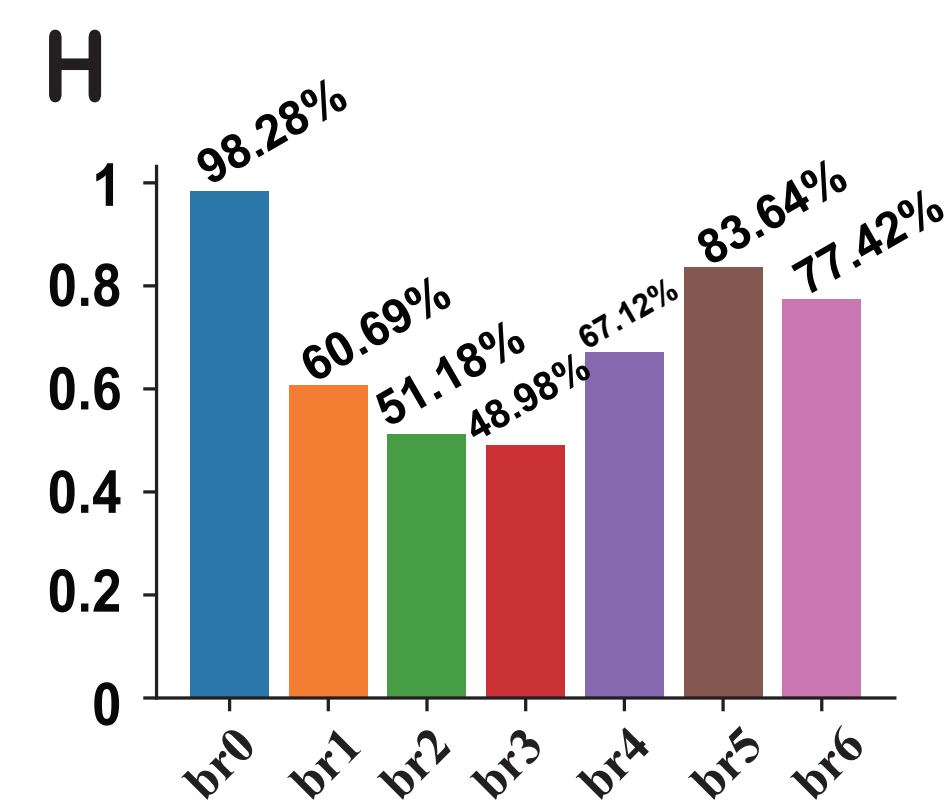
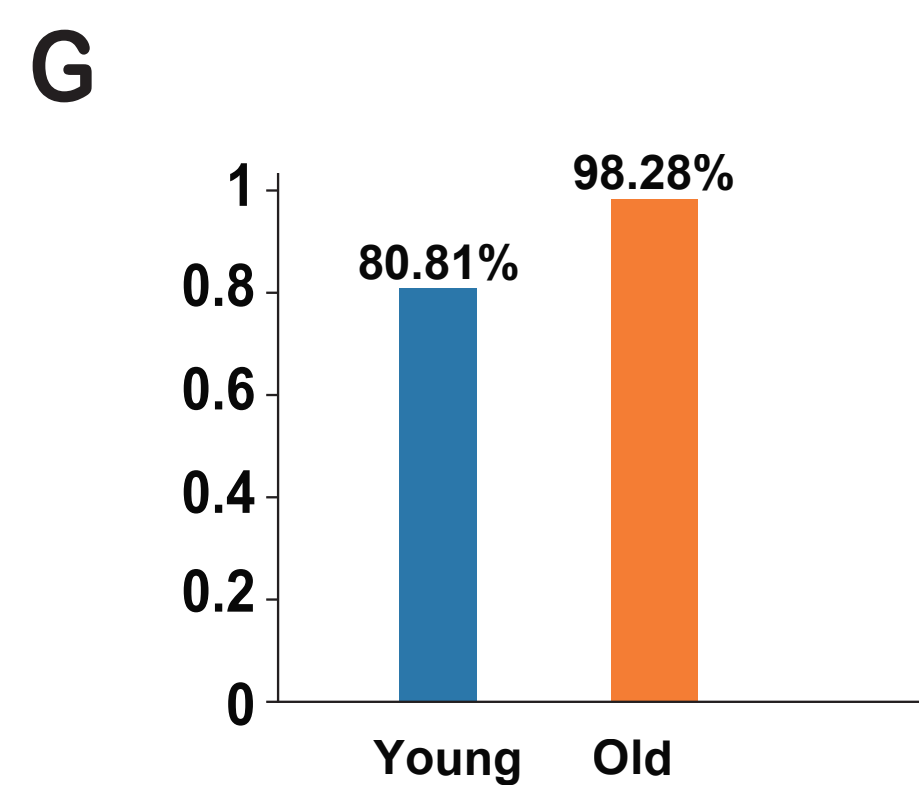
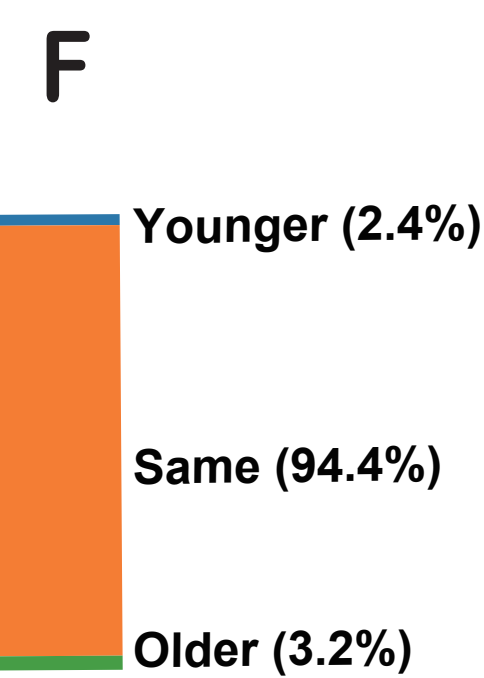
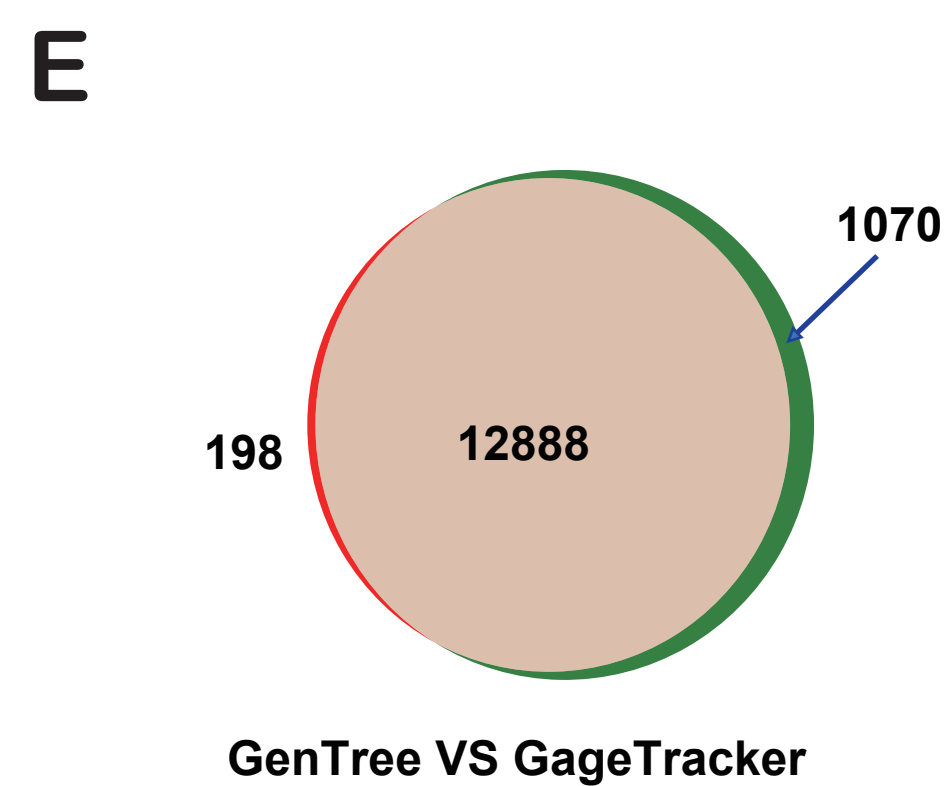
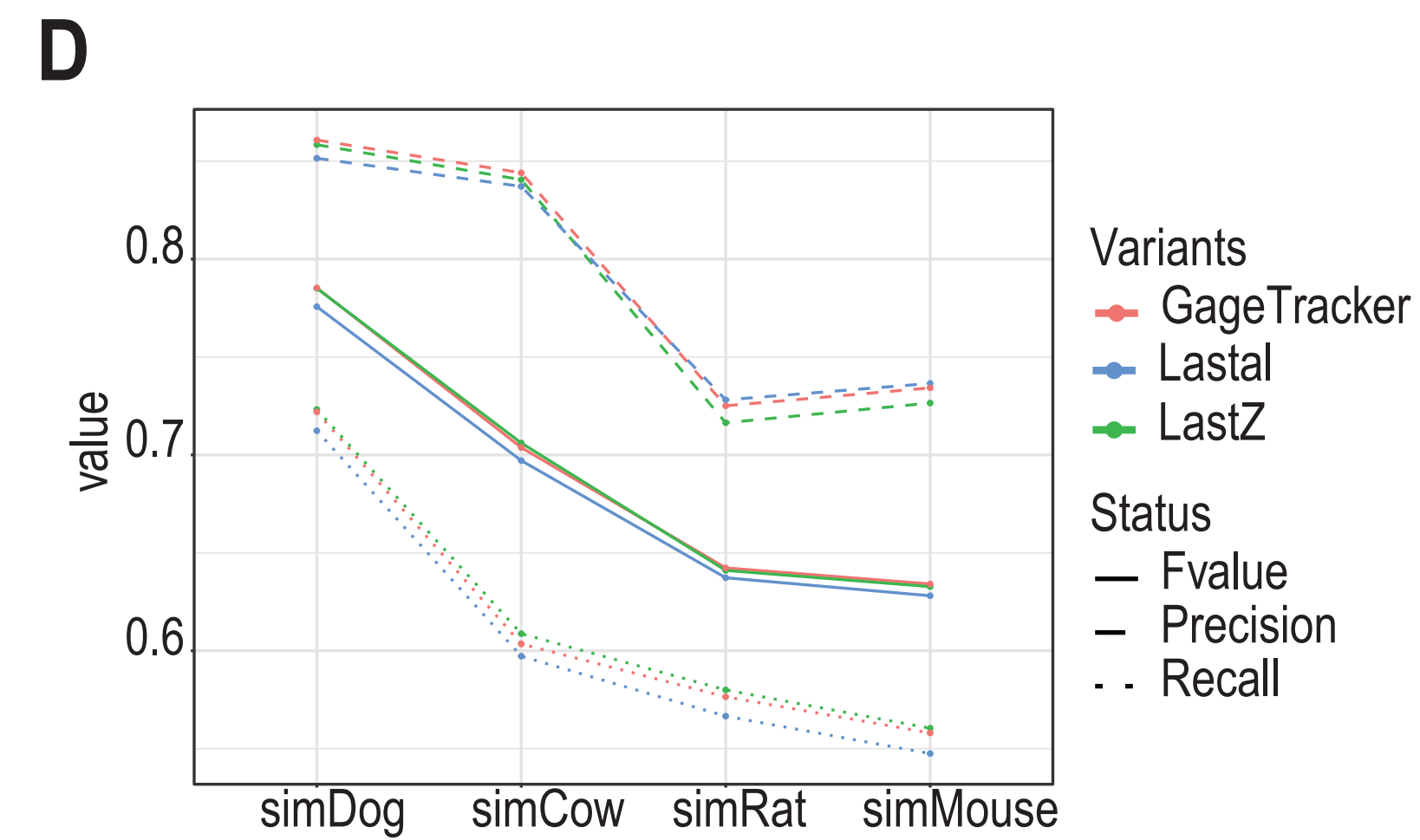
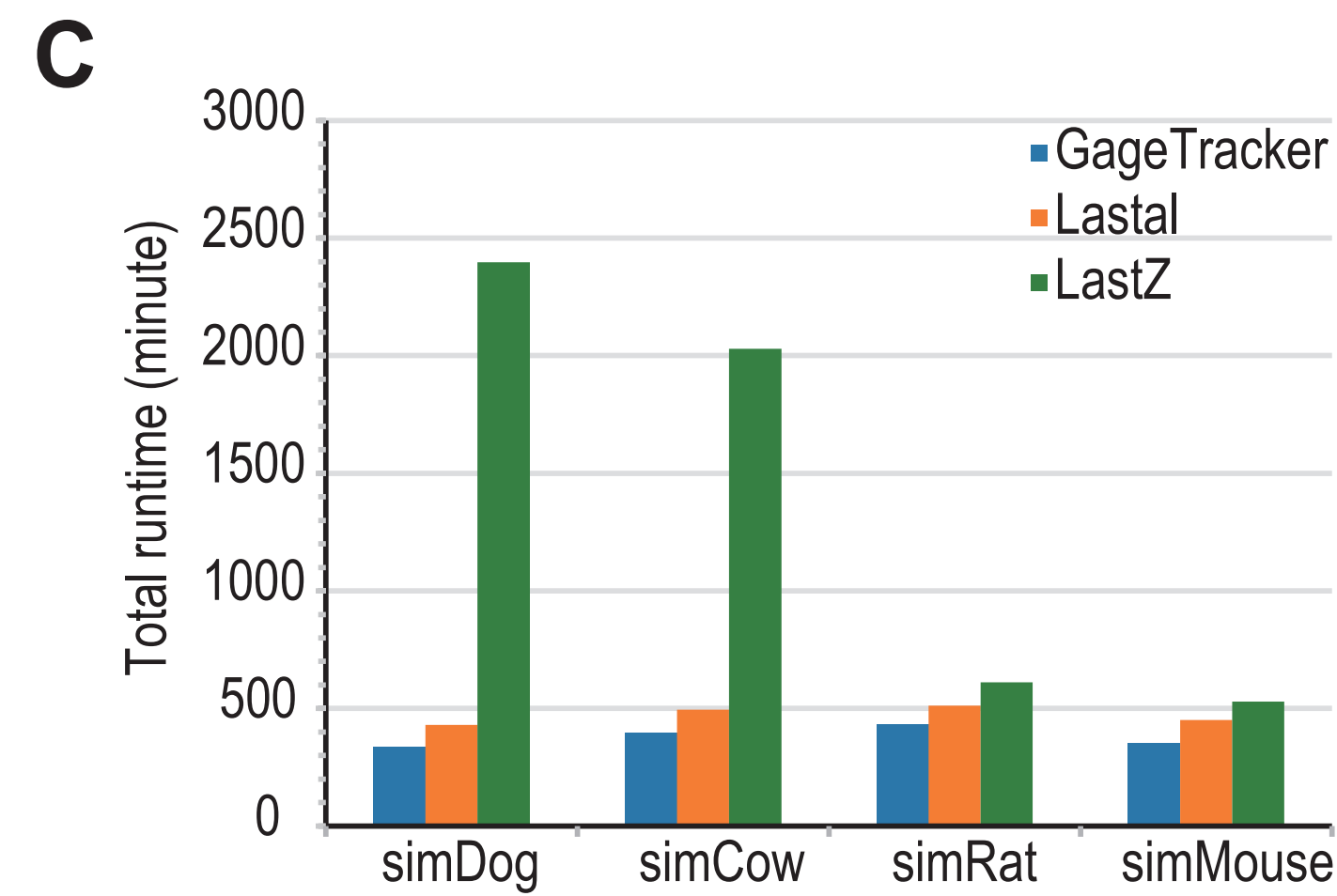
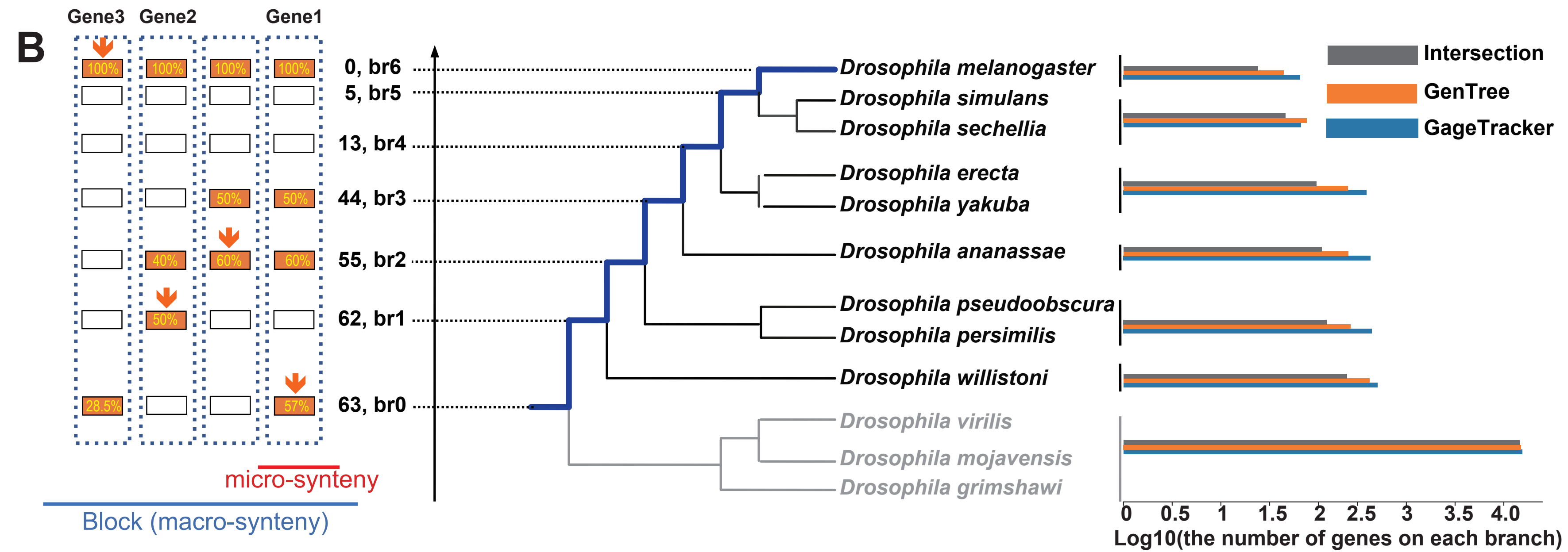
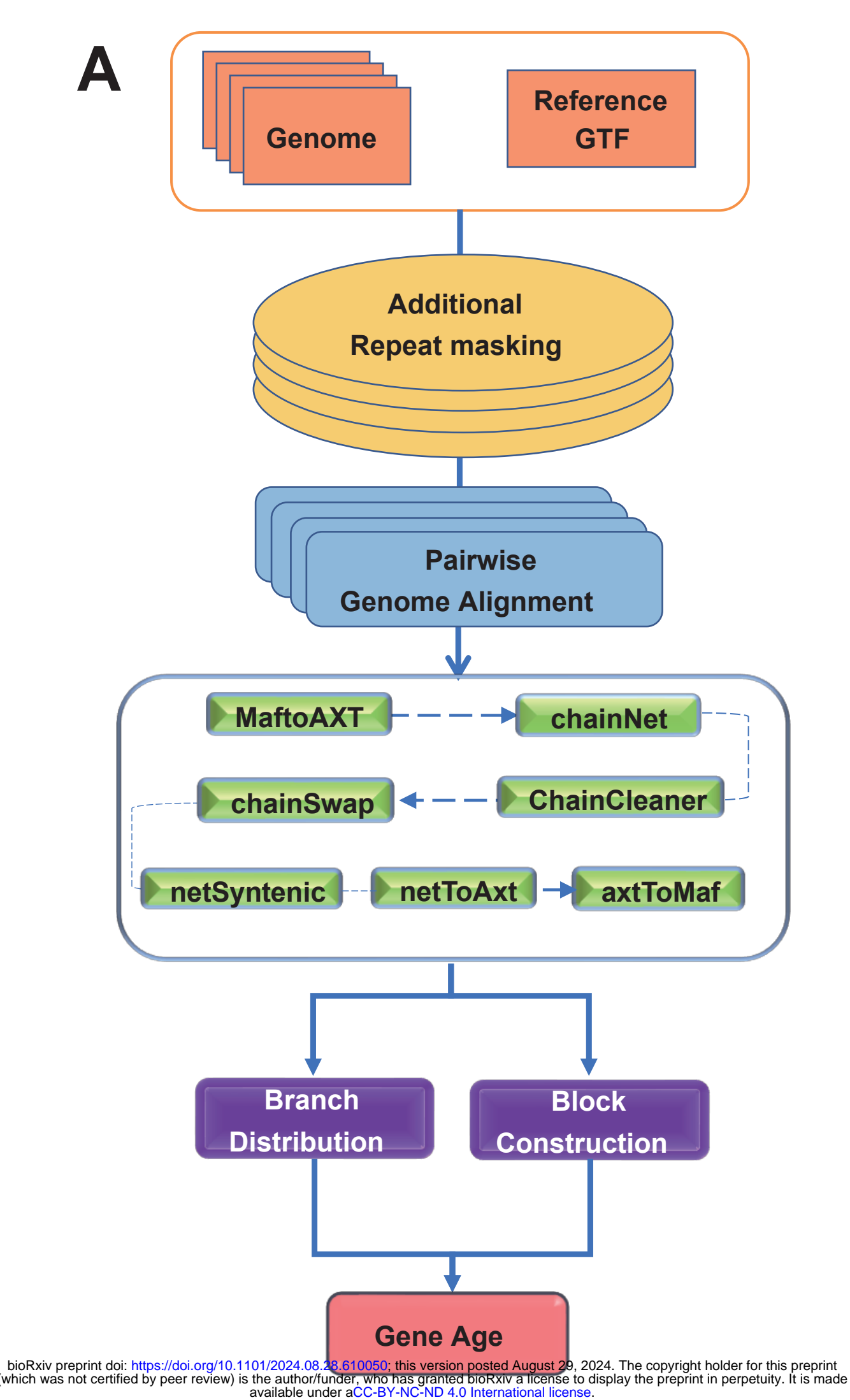
(last+without repeat-masking), GageTracker (last+repeat-masking/GageTracker -step1), and Lastz (LastZ + true repeats). LASTZ (LastZ+without repeat-masking) was not displayed in this figure due to its failure to finish the genome alignments in 7 days in the simulated data. We applied chains-and-nets and syntenic net algorithms to find reciprocal best alignments and Closure mafComparator to calculate precision, recall, and *F-score* values. (D) Precision, recall as well as *F-score* values of three aligners. (E) The comparison of Gene annotations employed by GenTree and GageTracker. (F) Comparison on the global level. (G) Comparison on the gene catalogs of young and old genes (nearly-detailed level). (H) Comparison on the gene catalogs of each branch (detailed level). (I) Cumulative curve of gene number with branch distance between the GageTracker and Gentree database.

Figure 2. Comparison of conflicting genes with branch distance between GageTracker and GenTree. (A) Cumulative curve of conflicting ages with branch distance between GageTracker vs. orthoDB and Gentree vs. orthoDB. *p-value* (0-2) represents the significance test in a shorter branch distance (0-2) group and *p-value* (3-6) represents the significance test in a larger branch distance (3-6) group. (B) Cumulative curve of conflicting ages with branch distance between GageTracker vs. Ensembl and Gentree vs. Ensembl. (C) Cumulative curve of conflicting ages with branch distance between GageTracker vs. FlyBase and Gentree vs. FlyBase. (D) The percentage of conflicting genes with multiple copies.

Figure 3. Expression pattern of conflicting genes in male of *D. melanogaster*. The tissue was marked in the top-left corner of each subfigure. The genes annotated as old gene in GageTracker ( branch 0 ) but young gene in GenTree ( from branches 1 to 6 ) formed Group A, while the genes annotated as young gene in GageTracker but old gene in GenTree formed Group B. The expression value was transformed by log function (log10). The red and green dash lines represent the mean expressions in group A and group B.

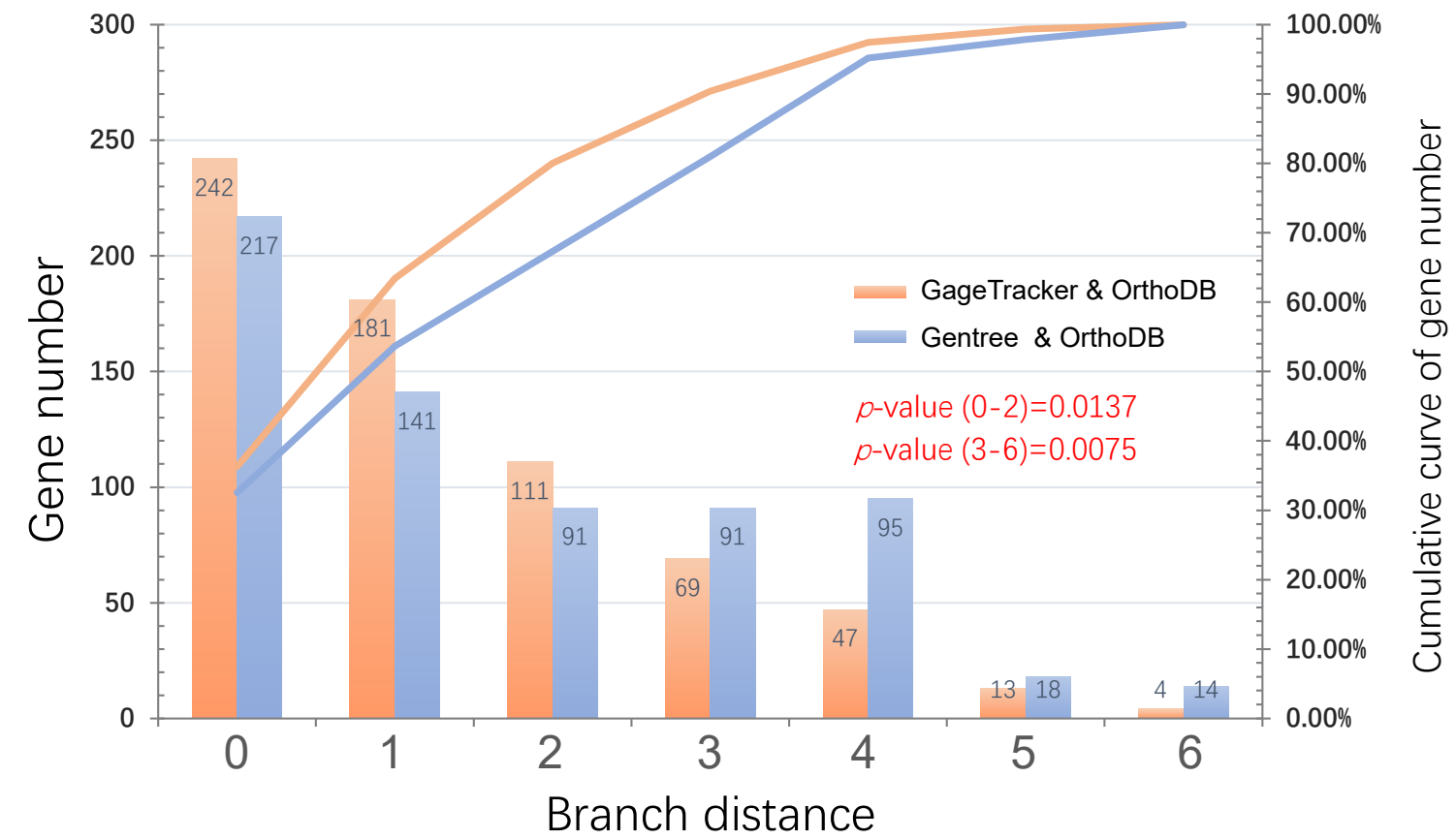
Figure 4. Expression pattern of conflicting genes in female of *D. melanogaster*. The tissue was marked in the top-left corner of each subfigure. The genes annotated as old gene in

GageTracker ( branch 0 ) but young gene in GenTree ( from branches 1 to 6 ) formed Group A, while the genes annotated as young gene in GageTracker but old gene in GenTree formed Group B. The expression value was transformed by log function ( $\log_{10}$ ). The red and green dash lines represent the mean expressions in group A and group B.

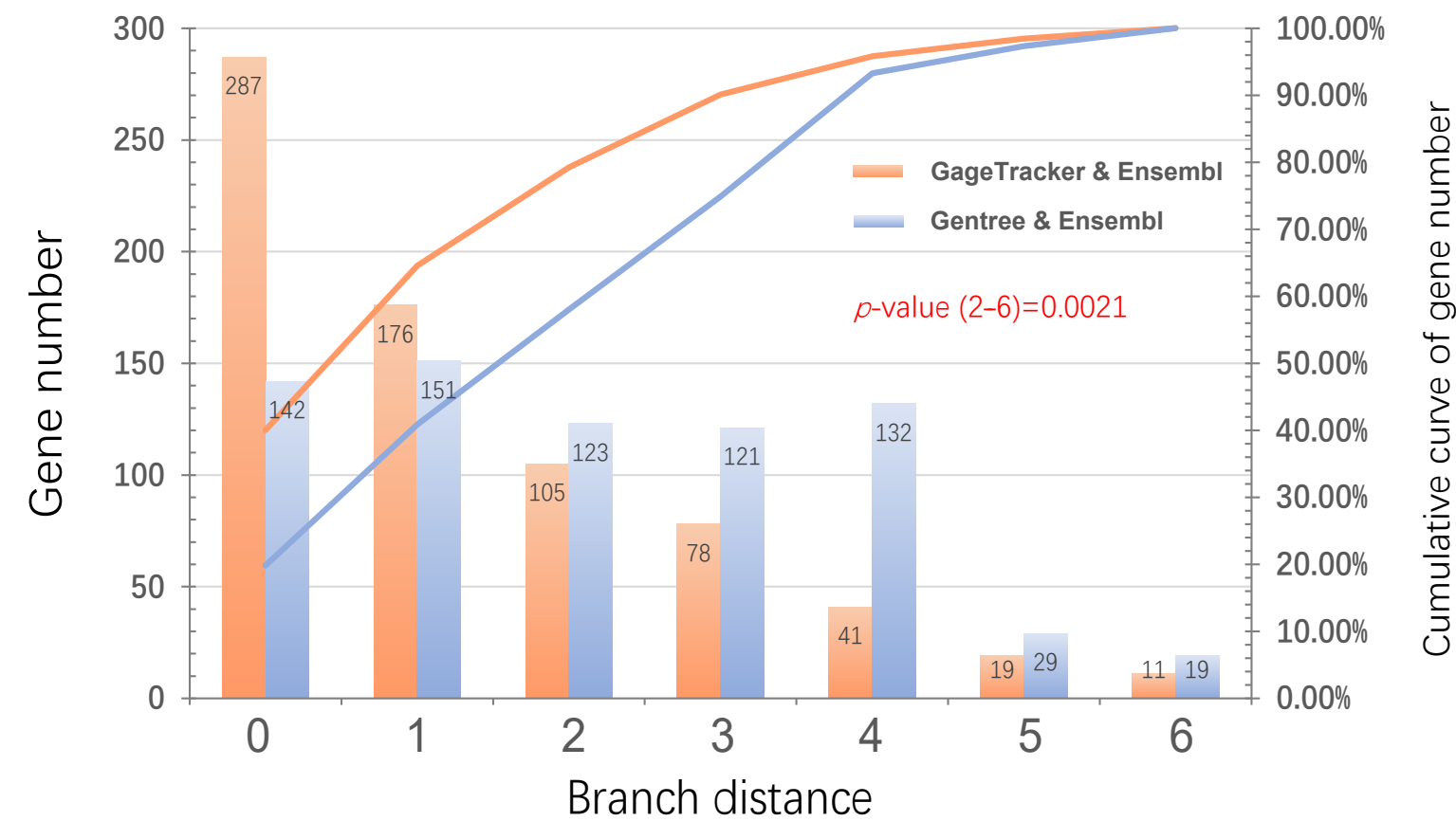


**A**

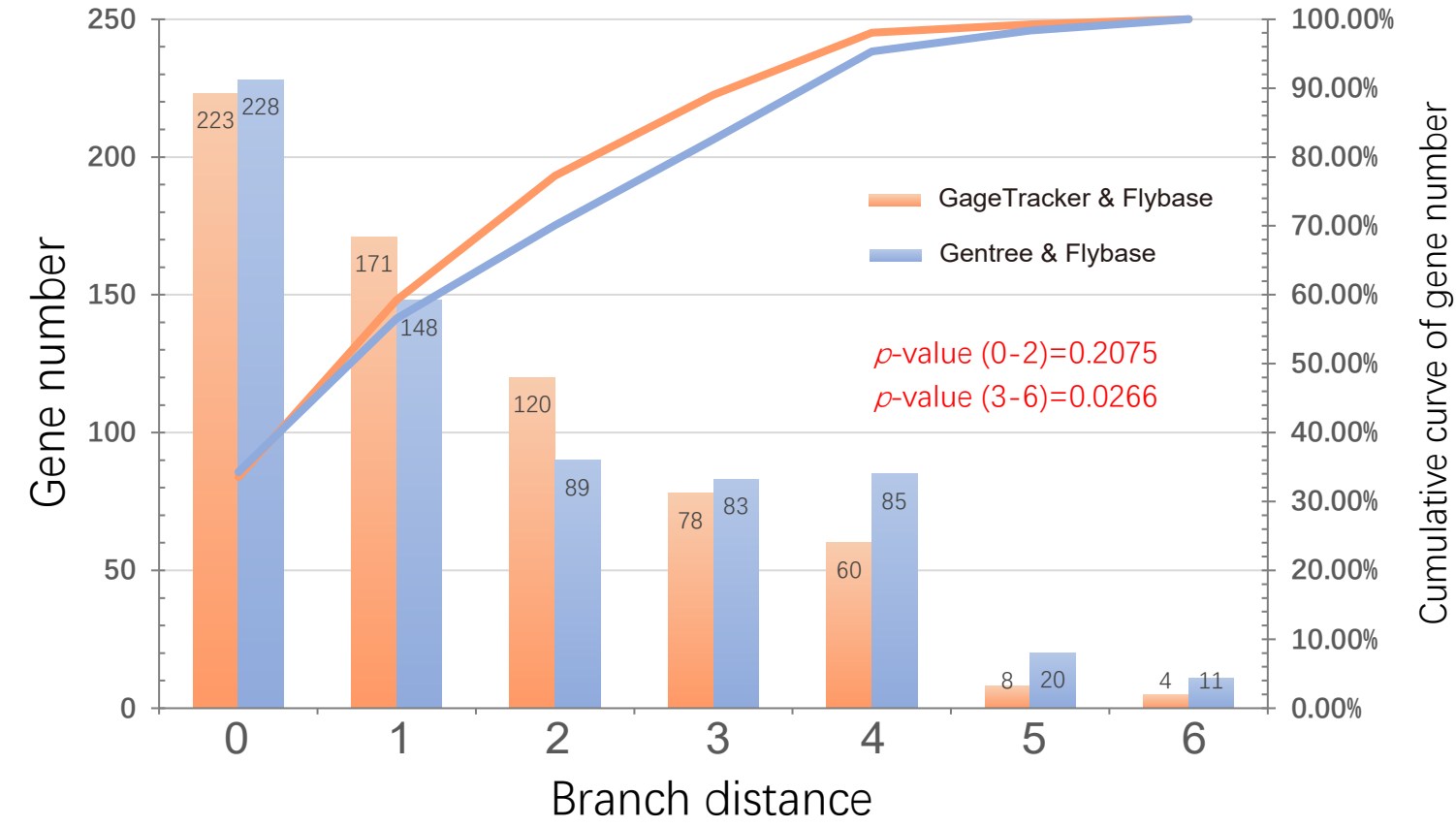
### Branch distance of conflicting genes with OrthoDB

**B**

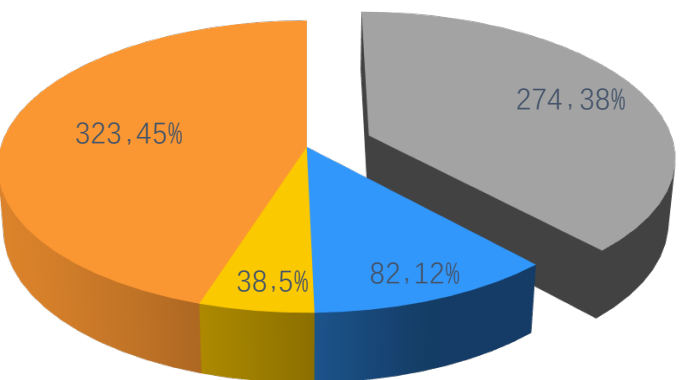
### Branch distance of conflicting genes with Ensembl

**C**

### Branch distance of conflicting genes with Flybase

**D**

Paralogous number of 717 conflicting genes



Paralogous number of 266 conflicting genes with distance >3

