

Supplementary Appendix A for *A new story of the four Hexapoda classes: Protura as the sister group to all other hexapods*

—Description of filtering strategy selection and matrix generation

Previous studies have shown that some gene properties are strongly correlated with phylogenetic signal<sup>1,2</sup>. Filtering genes based on their properties is a crucial step in phylogenomic studies<sup>2</sup>. In the last few years, much effort has been devoted to uncovering and understanding the roles of gene properties in phylogenomic reconstruction<sup>1–3</sup>. A total of ten properties were examined for each gene in this study, which we classified into two categories — four sequence-based (TAPER<sup>4</sup>, nRCFV<sup>5</sup>, SRH<sup>6</sup>, API<sup>7</sup>), and six tree-based properties (PSH<sup>8</sup>, ABS<sup>9</sup> values, DVMC<sup>10</sup>, treeness<sup>11</sup>, treeness over rcv<sup>11</sup>, and saturation<sup>12</sup>) — computed using custom scripts<sup>13</sup>. Their full descriptions are provided in Table 1). We tested the effect of these gene properties on three proxies for phylogenetic signal: Phylogenetic signal was evaluated based on ABS values, Robinson-Foulds (RF) distances<sup>14</sup>, and wASTRAL topologies.

First, filtering was performed with TAPER and TreeShink<sup>8</sup>, using default parameters (PSH analysis). Then, we focused on nRCFV and SRH that quantify compositional heterogeneity of alignments. We applied four thresholds (0.002, 0.003, 0.004, and 0.005) to nRCFV and examined the values for all genes, as shown in Figure S1a. Similarly, for assessing compositional homogeneity, we used two thresholds (0.05 and 0.1) for SRH and presented the values for all loci in Figure S1b. To evaluate the impact of these gene properties on phylogenetic reconstruction, we constructed wASTRAL trees using the matrices created from the sites retained by each threshold for nRCFV and SRH. Notably, all resulting trees strongly supported H1, as depicted in Figures S2 and S3. For detailed information on the number of loci, sites, and average locus length retained by each threshold, refer to Table 2. Additionally, the ABS values and RF distances of loci corresponding to different thresholds are displayed in Figures S4a, S4b, S5a, and S5b. To retain the maximum number of loci, we selected the threshold values of 0.004 for nRCFV (the "turning point" in Figure S1a) and 0.5 (the default value) for SRH. By combining these threshold values, we created the 'symtest-matrix' which included 949 loci, which served as the basis for subsequent analyses.

We then proceeded to investigate the six remaining properties of each analysed gene (API, ABS, DVMC, treeness, treeness over rcv, and saturation). Different thresholds were applied to filter the gene properties; the threshold and the resultant values are depicted in Figures S1c to S1h. For the API property, which was used to calculate the values of average pairwise identity (lower values indicating slowly evolving genes ) three thresholds were used: 0.6, 0.7, and 0.8. The wASTRAL trees constructed using the matrices created by the sites retained at the threshold lower than 0.6 strongly supported H3, while the other thresholds supported H2 (Figure S6). ABS filtering employed five thresholds: 55, 60, 65, 70, and 75. Loci of higher ABS values are thought to harbor more phylogenetic signal. The topology of the wASTRAL tree constructed using the matrices created by the sites retained at the threshold greater

than 55 supported H1, while the remaining thresholds supported H2 (Figure S7). We further quantified the genic deviation from the assumptions of a molecular clock (i.e., DVMC), under five thresholds: 0.5, 1, 1.5, 2, and 4 (Figure S1e). The wASTRAL trees strongly supported H1 (Figure S8). Treeness describes the signal-to-noise ratio in a phylogeny, whereby higher values of treeness are thought to be desirable<sup>11</sup> for phylogenetic inference. We set five thresholds: 0.1, 0.15, 0.17, 0.18, and 0.19 (Figure S1f). The wASTRAL trees constructed using the matrices created by the sites retained greater than each treeness threshold also strongly supported H1 (Figure S9). For treeness over rcv, three thresholds were used: 0.5, 0.8, and 1.1. When the threshold value greater than 1.1 was chosen, the topology of the wASTRAL tree constructed using this matrix created by the sites retained supported H2, whereas the other thresholds supported the topology [Collembola + Diplura] + [Protura + Insecta] (Figure S10). In the saturation analysis, two thresholds were considered: 0.1 and 0.4. Data with no saturation will have higher values. When the threshold > 0.1 was chosen, the topology of the wASTRAL tree constructed using the matrices created by the retained sites supported H1, while the threshold greater than 0.4 supported the topology of [Collembola + Diplura] + [Protura + Insecta] (Figure S11). All thresholds ensured that the number of retained loci was not less than 200, as specified in Table 2, which provides information on the number of loci, number of sites, and average locus length retained for each threshold. Additionally, ABS values and RF distances of the loci were calculated for each matrix (Figure S4, S5). Notably, significant differences were observed in ABS values (Figure S4d, S5d) and treeness (Figure S4f, S5f). However, for the API property, only ABS values showed significant differences (Figure S4c). Therefore, the ABS values and RF distance of loci do not differ significantly for different gene properties, even under different threshold (apart from API, ABS, and treeness).

In addition to the thresholds discussed above, the influence of missing data, i.e., the absence of some sites for some taxa, was considered. Five different thresholds (50, 60, 70, 80, and 90) were set, corresponding to the proportion of missing data in each dataset. Table 2 provides basic information on the thresholds, including the number of loci, number of sites, and average locus length retained for each threshold. The wASTRAL trees were constructed using the matrices created by the sites retained at each missing data threshold, and all results supported H1 (Figure S12). Notably, the ABS values and RF distances of loci did not differ significantly for the missing data gene properties, even with different thresholds (Figure S4h, S5h).

In summary, topological changes in the wASTRAL trees were observed for API and ABS at thresholds of 0.6 and 70, respectively. We concluded that ABS and API have an impact on phylogenetic inference in our dataset, but DVMC, and treeness, treeness over rcv, missing data, and saturation were not strongly correlated with phylogenetic signal. Two new matrices, namely 'API-matrix' and 'ABS-matrix', were created based on these thresholds. Importantly, we followed the rule of thumb that when selecting a threshold, the number of loci retained should not be less than half of the original dataset, which corresponded to retaining at least 475 loci in our case.

Table S1. Information on the ten gene properties (filtering strategies) used in this study.

| Property       | Name              | Description  |
|----------------|-------------------|--|
| Sequence-based | TAPER             | Two-dimensional Algorithm for Pinpointing ERrors   |
|                | nRCFV             | normalised Relative Compositional Frequency Variation  |
|                | API               | Average Pairwise Identity  |
|                | SRH               | Stationary, Reversible and Homogeneous   |
| Tree-based     | ABS               | Average Bipartition Support  |
|                | DVMC              | Degree of Violation of the Molecular Clock   |
|                | Treeness          | Proportion of sum of internal branch lengths over sum of all branch lengths across the maximum likelihood tree of a given alignment              |
|                | Treeness over RCV | Treeness divided by RCV  |
|                | PSH               | Potentially Spurious Homologs  |
|                | Saturation        | The sequences in multiple sequence alignments that have undergone numerous substitutions such that the distances between taxa are underestimated |

Table S2. Summary of the number of loci, number of sites, and average locus length for the matrices created using the sites retained at each threshold of nine distinct gene properties.

| Gene property     | threshold | Number of loci | Number of sites | Average locus length | wASTRAL topology |
|-------------------|-----------|----------------|-----------------|----------------------|------------------|
| nRCFV             | 0.002     | 589            | 290,160         | 492.63               | H1               |
|                   | 0.003     | 917            | 622,964         | 679.35               | H1               |
|                   | 0.004     | 994            | 742,413         | 746.89               | H1               |
|                   | 0.005     | 1,011          | 789,799         | 781.20               | H1               |
| SRH               | 0.05      | 949            | 685,769         | 722.62               | H1               |
|                   | 0.1       | 915            | 653,885         | 714.62               | H1               |
| API               | 0.6       | 475            | 374,376         | 788.16               | H3               |
|                   | 0.7       | 763            | 577,673         | 757.10               | H2               |
|                   | 0.8       | 904            | 664,304         | 734.84               | H2               |
| ABS               | 55        | 935            | 681,300         | 728.66               | H1               |
|                   | 60        | 885            | 664,515         | 750.86               | H2               |
|                   | 65        | 763            | 611,047         | 800.84               | H2               |
|                   | 70        | 554            | 499,622         | 901.84               | H2               |
|                   | 75        | 280            | 294,924         | 1,053.30             | H2               |
| DVMC              | 0.5       | 349            | 212,482         | 608.83               | H1               |
|                   | 1         | 749            | 509,478         | 680.21               | H1               |
|                   | 1.5       | 862            | 602,252         | 698.66               | H1               |
|                   | 2         | 900            | 637,395         | 708.21               | H1               |
|                   | 4         | 939            | 677,449         | 721.45               | H1               |
| treeness          | 0.1       | 947            | 684,669         | 722.98               | H1               |
|                   | 0.15      | 680            | 488,193         | 717.93               | H1               |
|                   | 0.17      | 452            | 329,502         | 728.98               | H1               |
|                   | 0.18      | 361            | 264,806         | 733.53               | H1               |
|                   | 0.19      | 270            | 200,448         | 742.40               | H1               |
| treeness over rcv | 0.5       | 883            | 604,438         | 684.52               | H1               |
|                   | 0.8       | 564            | 362,770         | 643.20               | H1               |
|                   | 1.1       | 249            | 151,690         | 609.19               | (C+D)+(P+I)      |
| missing data      | 50        | 948            | 685,373         | 722.97               | H1               |
|                   | 60        | 939            | 679,407         | 723.54               | H1               |
|                   | 70        | 900            | 635,031         | 705.59               | H1               |
|                   | 80        | 787            | 539,939         | 686.07               | H1               |
|                   | 90        | 470            | 333,303         | 709.16               | H1               |
| saturation        | 0.1       | 861            | 577,914         | 671.21               | H1               |
|                   | 0.4       | 371            | 195,529         | 527.03               | (C+D)+(P+I)      |

Note: H1, Collembola + (Protura + (Diplura + Insecta)); H2, (Collembola + Protura) + (Diplura + Insecta); H3, Protura + ((Collembola + Diplura) + Insecta).

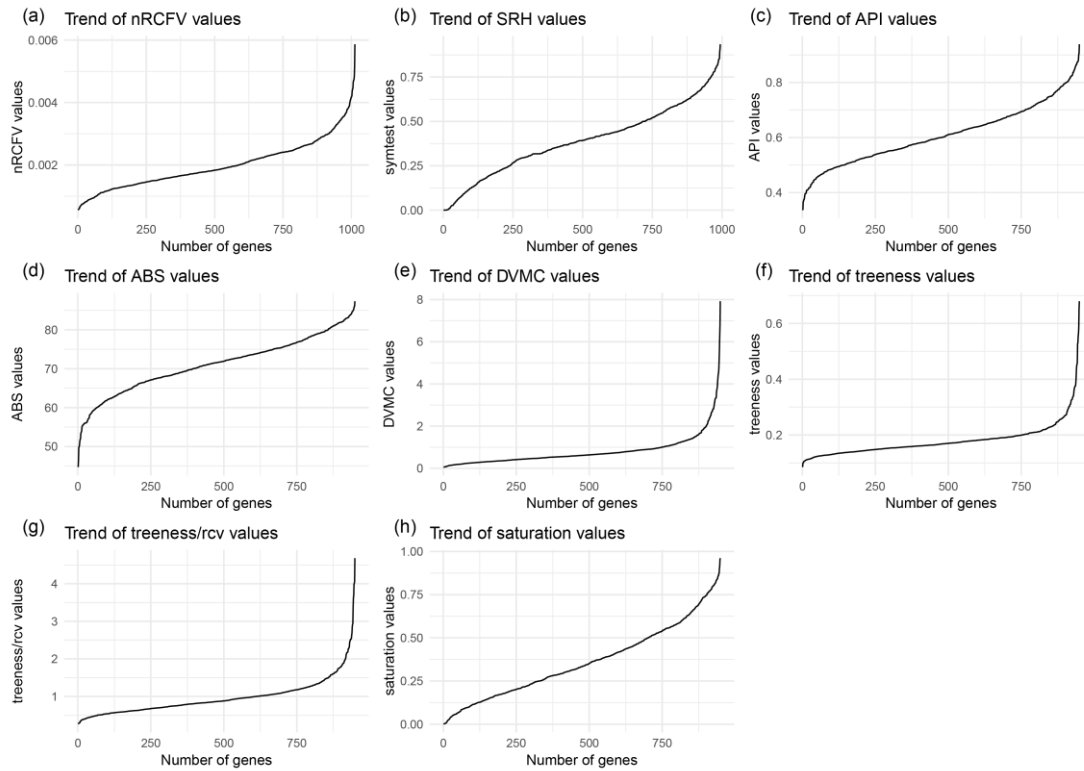


Figure S1. The values of different gene properties across different genes. (a) nRCFV (normalised Relative Compositional Frequency Variation); (b) SRH (Stationary, Reversible and Homogeneous); (c) API (Average Pairwise Identity); (d) ABS (Average Bipartition Support value); (e) DVMC (Degree of Violation of the Molecular Clock); (f) treeness; (g) treeness/rcv (treeness over rcv); (h) saturation.



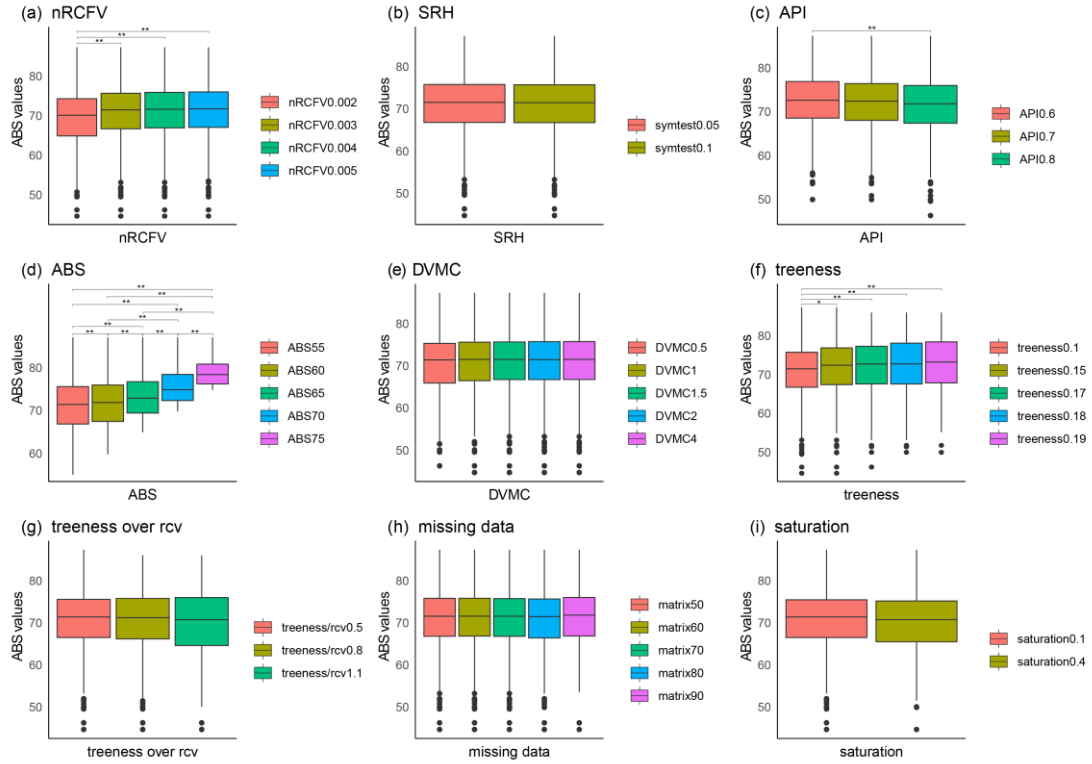


Figure S4. ABS values of matrix statistics created by loci under different threshold retention of (a) nRCFV, (b) SRH, (c) API, (d) ABS, (e) DVMC, (f) treeness, (g) treeness over rcv, (h) missing data, and (i) saturation. The top value corresponds to the p-value (\*  $p \leq 0.05$ ; \*\*  $p \leq 0.01$ ) show that there is significant difference between two matrices. No significant difference between two matrices is not shown.

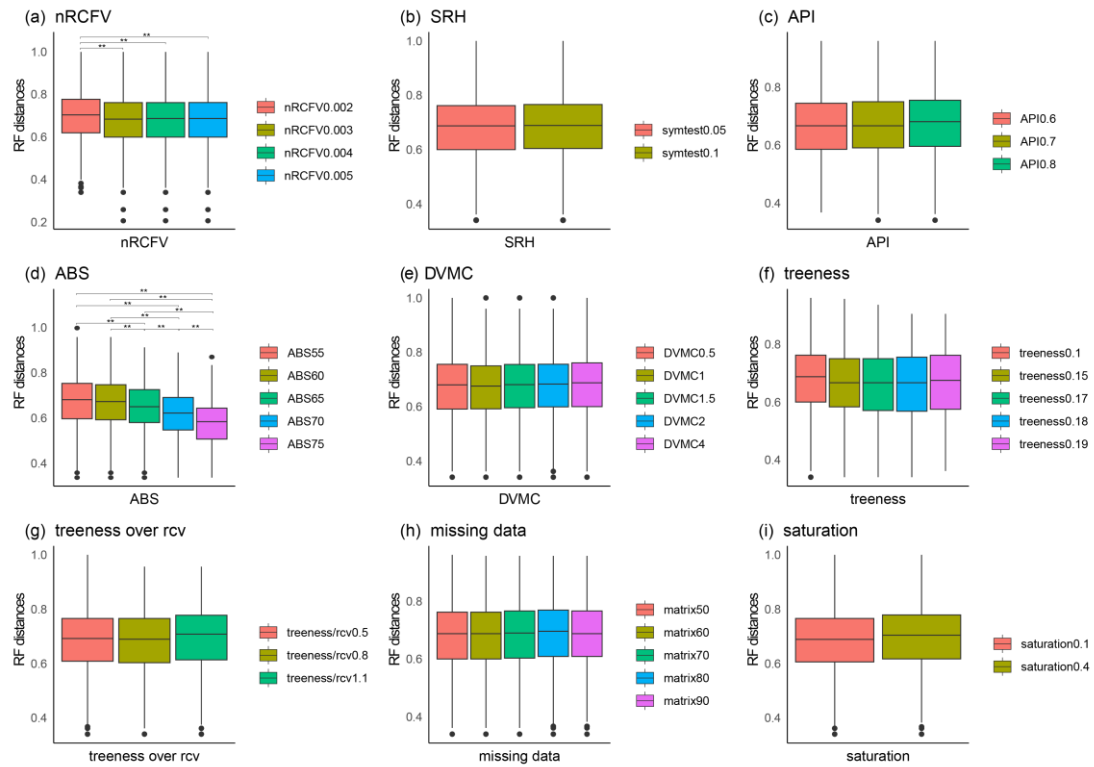


Figure S5. RF distances of matrix statistics created by loci under different threshold retention of (a) nRCFV, (b) SRH, (c) API, (d) ABS, (e) DVMC, (f) treeness, (g) treeness over rcv, (h) missing data, and (i) saturation. The top value corresponds to the p-value (\*  $p \leq 0.05$ ; \*\*  $p \leq 0.01$ ) show that there is significant difference between two matrices. No significant difference between two matrices is not shown.









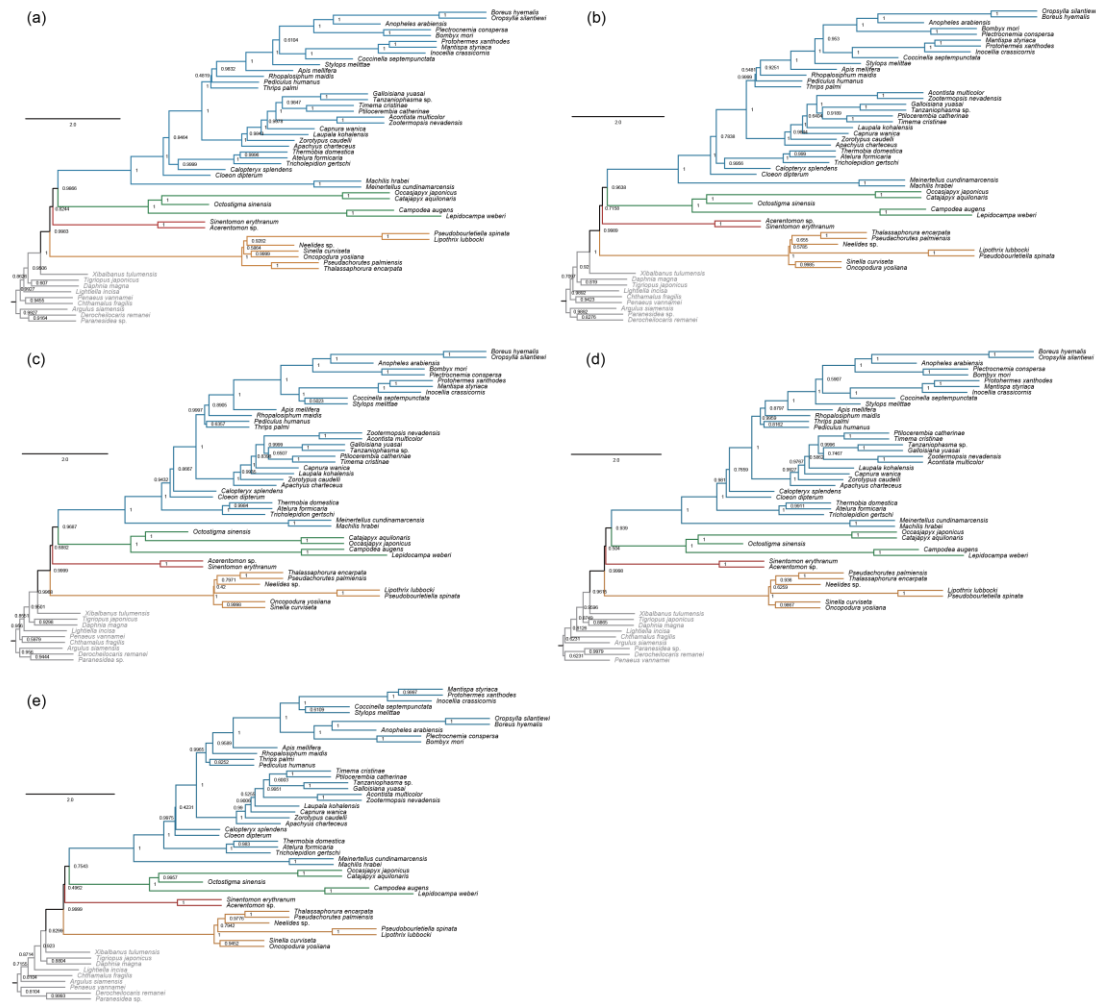
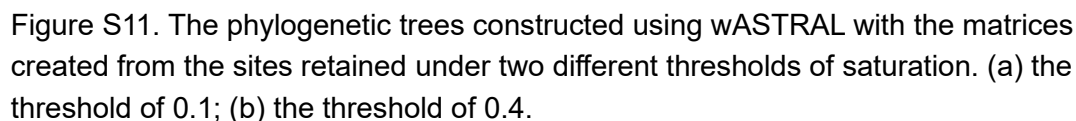
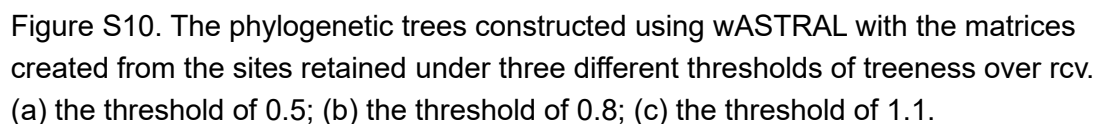


Figure S9. The phylogenetic trees constructed using wASTRAL with the matrices created from the sites retained under five different thresholds of treeness. (a) the threshold of 0.1; (b) the threshold of 0.15; (c) the threshold of 0.17; (d) the threshold of 0.18; (e) the threshold of 0.19.





## References

1. Kapli, P., Yang, Z., and Telford, M.J. (2020). Phylogenetic tree building in the genomic age. *Nat. Rev. Genet.* 21, 428–444. 10.1038/s41576-020-0233-0.
2. Young, A.D., and Gillung, J.P. Phylogenomics — principles, opportunities and pitfalls of big-data phylogenetics. *Syst. Entomol.* 45, 225–247. 10.1111/syen.12406.
3. Shen, X.-X., Steenwyk, J.L., and Rokas, A. (2021). Dissecting incongruence between concatenation- and quartet-based approaches in phylogenomic data. *Syst. Biol.* 70, 997–1014. 10.1093/sysbio/syab011.
4. Zhang, C., Zhao, Y., Braun, E.L., and Mirarab, S. (2021). TAPER: Pinpointing errors in multiple sequence alignments despite varying rates of evolution. *Methods Ecol. Evol.* 12, 2145–2158. 10.1111/2041-210X.13696.
5. Fleming, J.F., and Struck, T.H. (2023). nRCFV: a new, dataset-size-independent metric to quantify compositional heterogeneity in nucleotide and amino acid datasets. *BMC Bioinform.* 24, 145. 10.1186/s12859-023-05270-8.
6. Minh, B.Q., Schmidt, H.A., Chernomor, O., Schrempf, D., Woodhams, M.D., von Haeseler, A., and Lanfear, R. (2020). IQ-TREE 2: New models and efficient methods for phylogenetic inference in the genomic era. *Mol. Biol. Evol.* 37, 1530–1534. 10.1093/molbev/msaa015.
7. Chen, M.-Y., Liang, D., and Zhang, P. (2017). Phylogenomic resolution of the phylogeny of laurasiatherian mammals: exploring phylogenetic signals within coding and noncoding sequences. *Genome Biol. Evol.* 9, 1998–2012. 10.1093/gbe/evx147.
8. Mai, U., and Mirarab, S. (2018). TreeShrink: fast and accurate detection of outlier long branches in collections of phylogenetic trees. *BMC Genomics* 19, 272. 10.1186/s12864-018-4620-2.
9. Salichos, L., Rokas, A., and Rokas, A. (2013). Inferring ancient divergences requires genes with strong phylogenetic signals. *Nature* 497, 327–331. 10.1038/nature12130.
10. Liu, L., Zhang, J., Rheindt, F.E., Lei, F., Qu, Y., Wang, Y., Zhang, Y., Sullivan, C., Nie, W., Wang, J., et al. (2017). Genomic evidence reveals a radiation of placental mammals uninterrupted by the KPg boundary. *Proc. Natl. Acad. Sci.* 114, E7282–E7290. 10.1073/pnas.1616744114.
11. Phillips, M.J., and Penny, D. (2003). The root of the mammalian tree inferred from whole mitochondrial genomes. *Mol. Phylogenet. Evol.* 28, 171–185. 10.1016/S1055-7903(03)00057-5.
12. Philippe, H., Brinkmann, H., Lavrov, D.V., Littlewood, D.T.J., Manuel, M., Wörheide, G., and Baurain, D. (2011). Resolving difficult phylogenetic questions: why more sequences are not enough. *PLoS Biol.* 9, e1000602. 10.1371/journal.pbio.1000602.
13. Du, S., Ding, Y., Li, H., Zhang, A., Luo, A., Zhu, C., and Zhang, F. (2023). Construction of a phylogenetic matrix: scripts and guidelines for phylogenomics. *Zool. Syst.* 48, 107–116. 10.11865/zs.2023201.
14. Robinson, D.F., and Foulds, L.R. (1981). Comparison of phylogenetic trees. *Mat. Biosci.* 53, 131–147. 10.1016/0025-5564(81)90043-2.