

1 Title: **MATEdb2, a collection of high-quality metazoan proteomes across the Animal**
2 **Tree of Life to speed up phylogenomic studies**

3

4 Gemma I. Martínez-Redondo^{1*}, Carlos Vargas-Chávez¹, Klara Eleftheriadi¹, Lisandra
5 Benítez-Álvarez¹, Marçal Vázquez-Valls¹, Rosa Fernández^{1*}

6

7 ¹Metazoa Phylogenomics Lab. Biodiversity Program, Institute of Evolutionary Biology
8 (~~CSIC Spanish Research Council~~-University Pompeu Fabra). Passeig marítim de la
9 Barceloneta 37-49, 08003 Barcelona (Spain)

11 * Corresponding authors: gemma.martinez@ibe.upf-csic.es, rosa.fernandez@ibe.upf-csic.es

12

13

14 **Abstract**

15

16 Recent advances in high throughput sequencing have exponentially increased the number of
17 genomic data available for animals (Metazoa) in the last decades, with high-quality
18 chromosome-level genomes being published almost daily. Nevertheless, generating a new
19 genome is not an easy task due to the high cost of genome sequencing, the high complexity
20 of assembly, and the lack of standardized protocols for genome annotation. The lack of
21 consensus in the annotation and publication of genome files hinders research by making
22 researchers lose time in reformatting the files for their purposes but can also reduce the
23 quality of the genetic repertoire for an evolutionary study. Thus, the use of transcriptomes
24 obtained using the same pipeline as a proxy for the genetic content of species remains a
25 valuable resource that is easier to obtain, cheaper, and more comparable than genomes. In
26 a previous study, we presented the Metazoan Assemblies from Transcriptomic Ensembles
27 database (MATEdb), a repository of high-quality transcriptomic and genomic data for the two
28 most diverse animal phyla, Arthropoda and Mollusca. Here, we present the newest version
29 of MATEdb (MATEdb2) that overcomes some of the previous limitations of our database: (1)
30 we include data from all animal phyla where public data is available, (2) we provide gene
31 annotations ~~obtained~~ **extracted** from **the original GFF** genome **files** ~~obtained~~ using the same pipeline.
32 In total, we provide proteomes inferred from high-quality transcriptomic or genomic data for
33 almost 1000 animal species, including the longest isoforms, all isoforms, and functional
34 annotation based on sequence homology and protein language models, as well as the
35 embedding representations of the sequences. We believe this new version of MATEdb will
36 accelerate research on animal phylogenomics while saving thousands of hours of
37 computational work in a plea for open, greener, and collaborative science.

38 Introduction

39

40 In the midst of an explosion in the availability of genomic sequences, the
41 advancement of phylogenomic, phylotranscriptomic, and comparative genomic studies in
42 animals is hindered by the preprocessing and homogenization of the input data. With high-
43 quality chromosome-level genomes being published almost daily in the last few years, we
44 are gaining access to new biological knowledge that is helping to solve trickier scientific
45 questions, such as the identity of the sister taxon to *Bivalvia* (Song et al., 2023) or the
46 evolution of non-coding and repetitive regions (Osmanski et al., 2023). In addition, the use of
47 transcriptomes as a proxy of a species proteome continues to be a main source of proteome
48 data as a cheaper and easier alternative for phylogenetic inference (Erséus et al., 2020;
49 Mongiardino Koch et al., 2018; Zapata et al., 2014, among others) and gene repertoire
50 evolution (De Oliveira et al., 2016; Fernández & Gabaldón, 2020; Thoma et al., 2019) in
51 less-studied animals.

52

53 Together, these genomic and transcriptomic studies have provided a vast number of
54 resources for a plethora of animals that cannot be directly used in phylogenomic studies
55 before proper preprocessing. This is especially true for older datasets where data quality is
56 much lower and can have a high impact on the results obtained. Moreover, the use of
57 different computational pipelines for data processing makes data not comparable and prone
58 to false positives and negatives. For instance, the transcriptome assembly methodology
59 used can impact the comparability among different datasets (e.g. ~~in our experience for the~~
60 ~~subset of mollusk transcriptomes obtained from Krug et al. (2022)~~, the number of 'genes'
61 inferred with Trinity ~~is significantly different -p-value < 0.1 - to the ones we obtained~~, ~~Figure~~
62 ~~S1 may vary up to one order of magnitude depending on the version~~), while the 'ready-to-
63 use' protein files provided in some genome sequencing projects cannot be easily matched
64 with the other genome files for additional analyses due to different nomenclature across files.
65 This mainly impacts research groups with lower computational resources or experience who
66 cannot leverage the publicly available data into their research. To help alleviate these
67 issues, we previously published the Metazoan Assemblies from Transcriptomic Ensembles
68 database (MATEdb) containing high-quality transcriptome assemblies for 335 arthropods
69 and mollusks (Fernández et al., 2022). Here, we present its second version, MATEdb2, that
70 differs from the previous one in three main aspects: (1) we have increased the taxonomic
71 sampling to all animal phyla with high-quality data publicly available, and provide the first
72 transcriptomic sequences for some animal taxa; (2) we include a standardized pipeline for
73 obtaining the protein sequences from ~~GFF genomic files~~ instead of adding the
74 precomputed ~~protein files publicly available~~, making it easier to replicate ~~and combine with~~

75 | [the associated genomic sequence](#); (3) we provide the functional annotation of all proteins
76 | using a language-based new methodology that outperforms traditional methods_(Barrios-
77 | Núñez et al., 2024). We hope that this newer version of MATEdb accelerates research on
78 | animal evolution by providing a wider taxonomic resource of high-quality proteomes across
79 | the Animal Tree of Life.

80

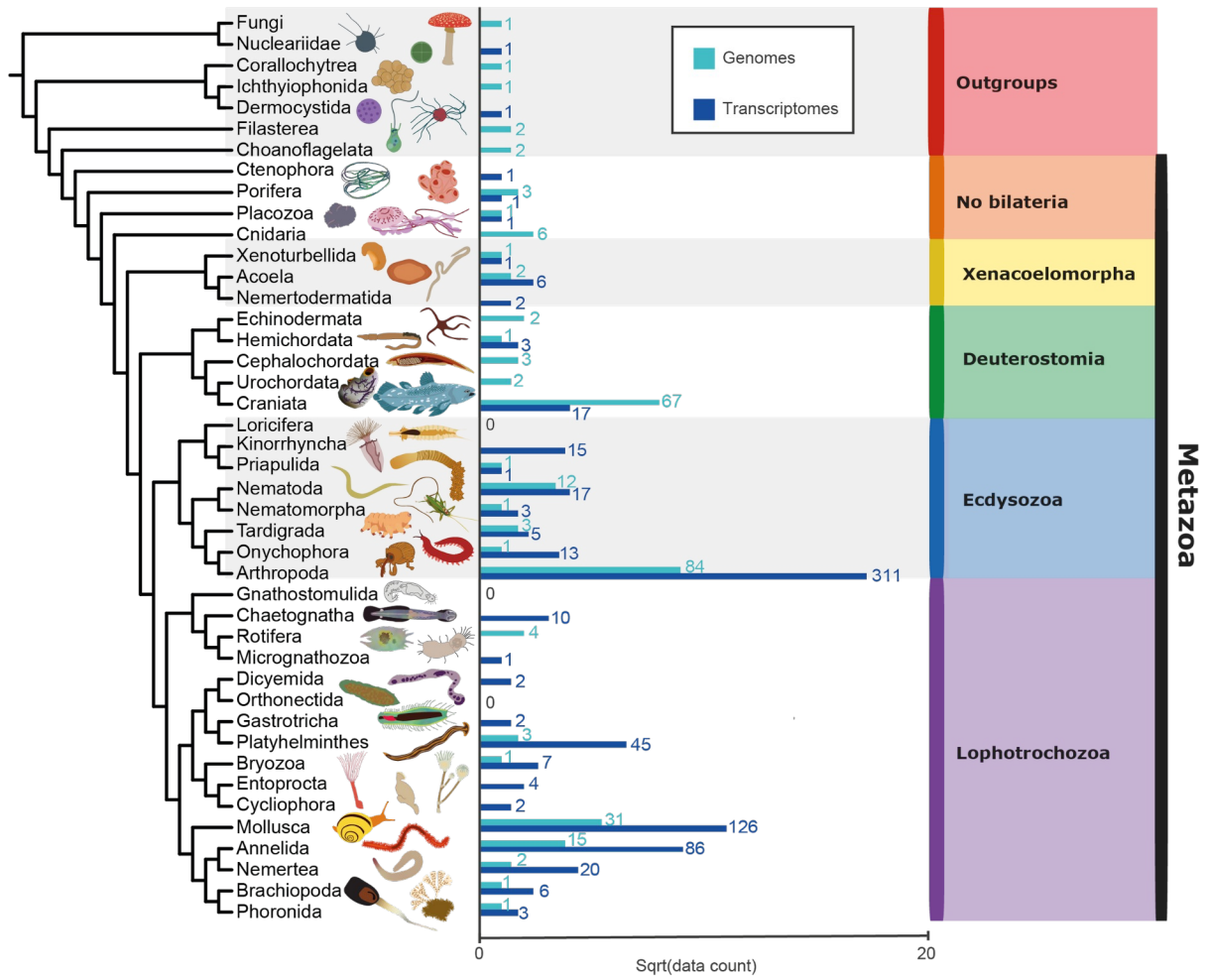
81 | **Material and Methods**

82

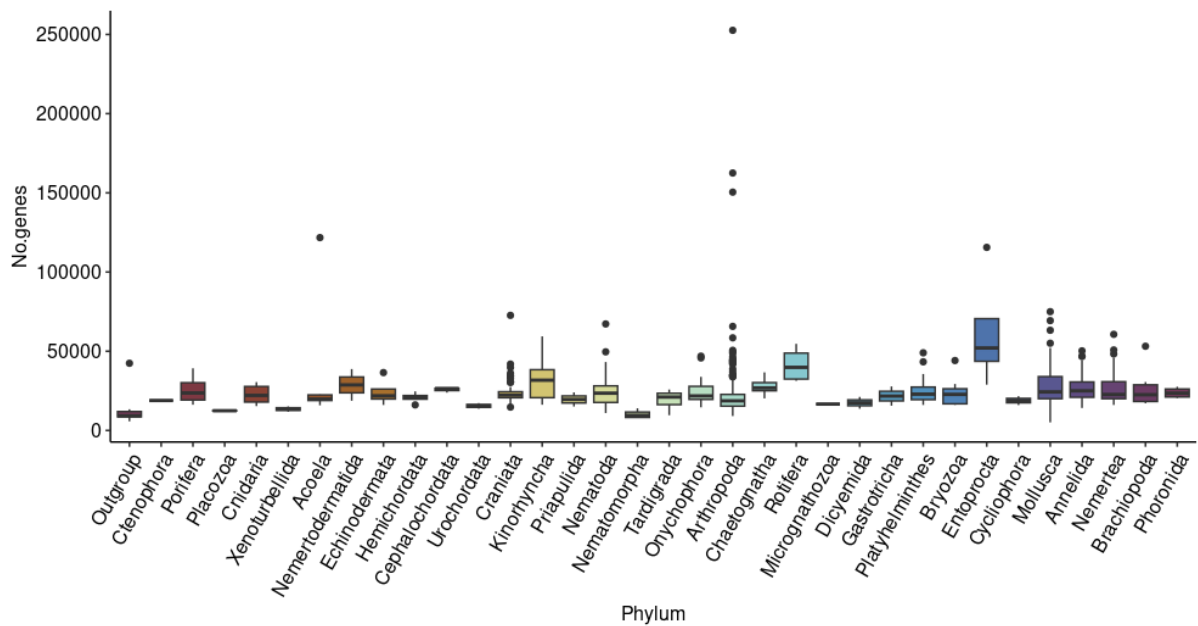
83 | *Increased taxonomic coverage*

84

85 | The first version of MATEdb (Fernández et al., 2022) included high-quality datasets
86 | from 335 species of arthropods and mollusks, with special attention to lineage representation
87 | within each phylum. Here, we provide a newer version of MATEdb that expands the
88 | taxonomic representation across the Animal Tree of Life by incorporating a total of 970
89 | species from virtually all animal phyla that have publicly available genomic or transcriptomic
90 | data, as well as some outgroup species relevant for understanding animal evolution. Taxon
91 | sampling tried to maximize the taxonomic representation within each phylum while
92 | considering the quality of the data. The number of proteomes per phylum included in
93 | MATEdb2 [and the distribution of gene number are](#) represented in Figures [1a and b](#), and
94 | the complete list of species and their metadata is included in Table S1.



97 **b.**



98

99 **Figure 1.** Taxonomic representation of species included in MATEdb2. **a.** The number of
100 datasets per phylum is separated by the type of data (genomes and transcriptomes in light
101 and dark blue, respectively). **b.** Distribution of the number of genes per proteome across
102 phyla.

103

104 *Improved analytical pipeline for genomes*

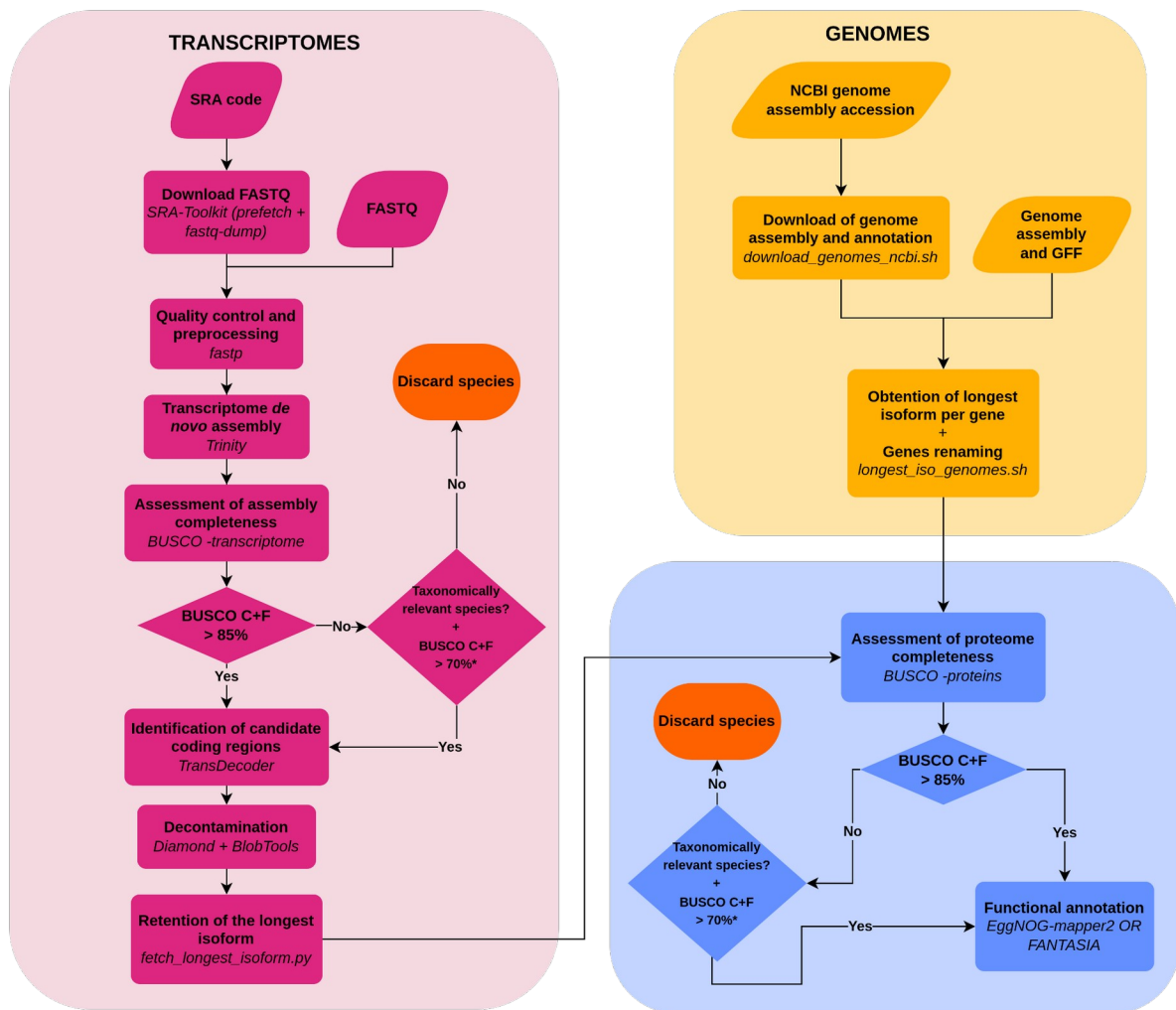
105

106 In the previous version of MATEdb (Fernández et al., 2022), we directly downloaded
107 the Coding DNA Sequences (CDS) and proteome files from the public repositories in the
108 case of genomes. However, a closer inspection of both files together with their
109 corresponding genome sequence and annotation revealed incongruences between them
110 that needed to be manually curated. Just looking at gene numbers, only 9 out of the 59
111 genomes we are keeping from the previous version of MATEdb had the same number of
112 protein-coding genes in the peptide and GFF files. In addition, most of the proteomes
113 differed in more than 1000 genes, with 15 having more than 10k of difference. This is caused
114 by the lack of consensus in the annotation and publication of genome files, with some
115 authors uploading modified versions of the protein sequences that do not map directly with

116 | the reported GFF and FASTA file, hindering the utility of those files for additional analyses.
117 | Another case we encountered when building the new version MATEdb2 was for the only
118 | chromosome-level ctenophore genome available back then (*Hormiphora californensis*), in
119 | which the proteins extracted directly from the GFF and FASTA file contained premature stop
120 | codons in virtually all the sequences, which made us discard this species. Moreover, even
121 | highly curated public databases can contain wrong or missing data, such as the case of *Apis*
122 | *mellifera* and *Anopheles gambiae*'s CDS file in Uniprot (UniProt Consortium, 2023)
123 | containing only a couple of sequences instead of the whole proteome. Therefore, we have
124 | included in the newer version of MATEdb a standardized pipeline for obtaining the CDS and
125 | protein files using directly the FASTA and GFF files of the corresponding genome.

126

127 | The analytical pipeline of MATEdb2 is shown in Figure 2. In brief, the differences with
128 | the pipeline depicted in MATEdb (Fernández et al., 2022) are the following: (1) we included
129 | a standardized pipeline for obtaining the longest isoform from genomes; (2) for a few
130 | exceptions, we lowered the threshold used to consider a dataset as high-quality to 70% C+F
131 | (complete plus fragmented) BUSCO score (Manni et al., 2021), as the original 85%
132 | threshold was too restrictive when prioritizing a wide taxonomic sampling and the inclusion
133 | of biologically interesting species that are not widely studied. Further details about the
134 | pipeline are shown in Figure 2.



135

136

137 **Figure 2.** Pipeline followed to generate the MATEdb2 database. All steps differing from the
 138 MATEdb original pipeline are discussed in detail in the main text.

139

140 *Compilation of genomic data*

141

142 Genome assembly (FASTA) and annotation (GFF) files for each species were
 143 downloaded through NCBI Datasets (Sayers et al., 2023) or from the direct URL download
 144 link for genomes available in other repositories. The database source for each species is
 145 referenced in the Supplementary Material, Table S1, while the bash script
 146 “download_genomes.sh” used to automatize the downloading of several files is included in
 147 the GitHub repository and the Singularity container (See Data Availability).

148

149 Once downloaded, we used AGAT ([Dainat et al., 2023](#) ~~Creators Jacques Dainat1~~
 150 ~~Darío Hereñú~~ ~~Dr. K. D. Murray2~~ ~~Ed Davis3~~ ~~Kathryn Crouch4~~ ~~LucileSol~~ ~~Nuno Agostinho5~~
 151 ~~pascal-git~~ ~~Zachary Zollman~~ ~~tayyrov~~ ~~Show affiliations~~ 1. IRD 2. Max Planck Institut für

152 | ~~Entwicklungsbiologie, Tübingen 3. Oregon State University 4. @VEuPathDB 5. European~~
153 | ~~Bioinformatics Institute | EMBL-EBI, n.d.~~) to obtain the GFF containing only the longest
154 | isoforms which was then used to get the FASTA file with the longest protein sequence for
155 | each gene (and its corresponding CDS). In addition, we renamed the sequences to match
156 | the structure used in the transcriptomic part of the MATEdb2 pipeline and obtained a
157 | conversion file to keep track of the original names. These steps were performed using a
158 | custom bash script “longest_iso_genomes.sh”, also included in the GitHub repository and
159 | container.

160

161 | Finally, gene completeness was assessed using BUSCO in protein mode against the
162 | metazoa_odb10 reference set (except for the outgroup species, where eukaryota_odb10
163 | was used). More than 75% of our species passed the threshold of 85% complete plus
164 | fragmented used in MATEdb (Fernández et al., 2022). The remaining 25% includes almost
165 | all representatives of tardigrades, annelids, nematodes, acoels, and some representatives of
166 | other phyla (see Table S1). As we want to maximize the taxon representation of animal
167 | lineages while keeping datasets of high quality, we lowered the threshold value to 70% in
168 | these cases, a value ~~that has been~~ previously used in other studies, ~~I, as t~~ these values may
169 | ~~indeed~~ represent biological features of the genomes of these lineages (Barreira et al., 2021)
170 | ~~or just a lack of representation of the lineage in the BUSCO reference datasets~~. As an
171 | exception, after this new threshold, 8 animal and 2 outgroup transcriptome assemblies have
172 | been included with a slightly lower BUSCO score due to their taxonomic relevance (e.g.,
173 | they were ~~one of~~ the only ~~two~~ representatives of their lineage, such as in the case of the
174 | ~~Priapulida hagfish~~). ~~A list of discarded datasets can also be found in Table S2.~~

175

176 | ***Functional annotation of the gene repertoire***

177

178 | The longest isoform gene list for each dataset was annotated with the homology-
179 | based software eggNOG-mapper v2 (Cantalapiedra et al., 2021) and the FANTASIA pipeline
180 | (<https://github.com/MetazoaPhylogenomicsLab/FANTASIA>). FANTASIA is a pipeline that
181 | allows the annotation of whole proteomes using GOPredSim (Littmann et al., 2021), a
182 | protein language-based method that transfers GO terms based on embedding similarity. In
183 | brief, embeddings are vectorized representations of protein sequences generated using
184 | protein language models, such as ProtT5 (Elnaggar et al., 2022), that consider protein
185 | sequences as sentences and apply natural language processing tools to extract information
186 | from them. Here, besides the GO terms predicted by FANTASIA, we provide the raw per-
187 | protein ProtT5 embeddings. More details about the pipeline, the method or the

188 benchmarking and comparison with homology-based methods can be checked [elsewhere](#)
189 (Barrios-Núñez et al., 2024;) and Martínez-Redondo et al. 2024).

190

191 **Database availability**

192

193 *Scripts and commands*

194

195 Scripts and commands in the pipeline and the supplementary [data](#) (Tables S1-2,
196 [Figure S1](#)) can be found in the following repository:

197 <https://github.com/MetazoaPhylogenomicsLab/MATEdb2>

198

199 *Files deposited in the repository*

200

201 For transcriptomes, the data repository contains (1) de novo transcriptome
202 assemblies, (2) their candidate coding regions within transcripts (both at the level of
203 nucleotide and amino acid sequences), (3) the coding regions filtered using their
204 contamination profile (ie, only metazoan content or eukaryote for outgroups), (4) the longest
205 isoforms of the amino acid candidate coding regions, (5) the gene content completeness
206 score as assessed against the BUSCO reference sets, and (6) orthology and protein
207 language-based gene annotations, and per-protein ProtT5 embeddings. In the case of
208 genomes, only files (4), (5), and (6) are provided in MATEdb2, together with a filtered
209 version of the file (3) with just the longest CDS per gene.

210

211 [The database is hosted on our own server and will be there indefinitely. The
212 database will be expanded as we incorporate new datasets from underrepresented lineages,
213 such as nematodes, or as requested to be incorporated by the scientific community if
214 resources allow it. Links for downloading can be found in the following file in the GitHub
215 repository:](#)

216 <https://github.com/MetazoaPhylogenomicsLab/MATEdb2/blob/main/linksforMATEdb2.txt>

217

218 *Software availability*

219

220 We provide a Singularity container for easy implementation of the tools used to generate the
221 files in the database with the appropriate software versions along with their dependencies
222 (<https://cloud.sylabs.io/library/klarael.metazomics/matedb2/matedb2.sif>). The software
223 included is the following: SRA Toolkit v2.10.7 (<http://ncbi.github.io/sra-tools/>), fastp v0.20.1
224 (<https://github.com/OpenGene/fastp>; Chen et al., 2018), Trinity v2.11.0 ([Grabherr et al.](#)

225 [2011](#)), BUSCO v5.3.2 ([Manni et al., 2021](#)), TransDecoder v5.5.0
226 (<https://github.com/TransDecoder/TransDecoder>), Diamond v2.0.8 ([Buchfink, Xie, & Huson,](#)
227 [2015](#)), BlobTools v2.3.3 ([Challis et al., 2020](#)), NCBI datasets v13.42.0, eggNOG-mapper
228 v2.1.9 ([Cantalapiedra et al., 2021](#)), seqkit v2.1.0 ([Shen, Sipos, and Zhao, 2024](#)), AGAT
229 v0.9.1 ([Daniat et al., 2023](#)), as well as some custom scripts.

230

231

232 [Discussion](#)

233

234 [We presented here the second version of MATEdb \(MATEdb2\), with almost 1000 animal](#)
235 [species data. This newer version overcomes some of the previous restrictions of our](#)
236 [database, including the restricted taxonomic representation of only arthropods and mollusks,](#)
237 [and the use of previously preprocessed peptide and CDS files for genomes. Nevertheless, it](#)
238 [is not devoid of limitations. The main limitation of this newer version are still the genome](#)
239 [annotations. Even though we use an alternative approach that considers the](#)
240 [incongruences found in some of the publicly preprocessed files, there may still be](#)
241 [biases between the proteomes. These biases are caused by the heterogeneity of](#)
242 [genome annotation methodologies, which can affect downstream analyses, such as](#)
243 [ortholog inference \(Weisman et al., 2022\). These biases are typically ignored by](#)
244 [phylogenomic studies that use publicly available preprocessed files. Nevertheless,](#)
245 [correcting for this limitation by re-annotating the genomes using the same](#)
246 [methodology is computationally expensive and is still biased toward species where](#)
247 [additional data that improves this annotation \(e.g. RNA-seq\) is available.](#)

248

249

250 **Author contributions**

251

252 This database results from the collaborative effort of lab members from the Metazoa
253 Phylogenomics Lab to offer the scientific community the possibility to reuse some of the data
254 generated for their projects. GIMR, CVC, KE, LBA, and MVV contributed assemblies to the
255 data repository. GIMR created the pipeline custom scripts for the genome data analyses and
256 designed the MATEdb logo. KE created the Singularity container. CVC and RF contributed
257 to the creation and management of the database. CVC created and curated the Github
258 repository. RF provided resources and supervised the project. GIMR wrote the first version
259 of the manuscript. All authors revised and approved the final version of the manuscript.

260

261

262 Acknowledgments

263

264 GIMR acknowledges the support of Secretaria d'Universitats i Recerca del Departament
265 d'Empresa i Coneixement de la Generalitat de Catalunya and ESF Investing in your future
266 (grant 2021 FI_B 00476). RF acknowledges support from the following sources of funding:
267 Ramón y Cajal fellowship (grant agreement no. RYC2017-22492 funded by MCIN/AEI
268 /10.13039/501100011033 and ESF 'Investing in your future'), the Agencia Estatal de
269 Investigación (project PID2019-108824GA-I00 funded by
270 MCIN/AEI/10.13039/501100011033), the European Research Council (this project has
271 received funding from the European Research Council (ERC) under the European's Union's
272 Horizon 2020 research and innovation programme (grant agreement no. 948281)), the
273 Human Frontier Science Program (grant no. RGY0056/2022) and the Secretaria
274 d'Universitats i Recerca del Departament d'Economia i Coneixement de la Generalitat de
275 Catalunya (AGAUR 2021-SGR00420). We also thank Centro de Supercomputación de
276 Galicia (CESGA) and the HPC Drago from the Centro Superior de Investigaciones
277 Científicas for access to computer resources.

278

279

280 References

281 Barreira, S. N., Nguyen, A.-D., Fredriksen, M. T., Wolfsberg, T. G., Moreland, R. T., &

282 Baxevanis, A. D. (2021). AniProtDB: A Collection of Consistently Generated Metazoan

283 Proteomes for Comparative Genomics Studies. *Molecular Biology and Evolution*,

284 38(10), 4628–4633. <https://doi.org/10.1093/molbev/msab165>

285 Barrios-Núñez, I., Martínez-Redondo, G. I., Medina-Burgos, P., Cases, I., Fernández, R., &

286 Rojas, A. M. (2024). Decoding proteome functional information in model organisms

287 using protein language models. In *bioRxiv* (p. 2024.02.14.580341).

288 <https://doi.org/10.1101/2024.02.14.580341>

289 [Buchfink, Benjamin, Chao Xie, and Daniel H. Huson. 2015. "Fast and Sensitive Protein](#)

290 [Alignment Using DIAMOND." *Nature Methods* 12 \(1\): 59-60.](#)

291 <https://doi.org/10.1038/nmeth.3176>

292 Cantalapiedra, C. P., Hernández-Plaza, A., Letunic, I., Bork, P., & Huerta-Cepas, J. (2021).

293 eggNOG-mapper v2: Functional Annotation, Orthology Assignments, and Domain

294 Prediction at the Metagenomic Scale. *Molecular Biology and Evolution*, 38(12), 5825–
295 5829. <https://doi.org/10.1093/molbev/msab293>

296 [Challis, Richard, Edward Richards, Jeena Rajan, Guy Cochrane, and Mark Blaxter. 2020.](#)
297 [“BlobToolKit - Interactive Quality Assessment of Genome Assemblies.” G3 10 \(4\): 1361-](#)
298 [74. <https://doi.org/10.1534/g3.119.400908>](#)

299 [Chen, S., Zhou, Y., Chen, Y., & Gu, J. \(2018\). fastp: an ultra-fast all-in-one FASTQ](#)
300 [preprocessor. *Bioinformatics*, 34\(17\), i884–i890.](#)
301 <https://doi.org/10.1093/bioinformatics/bty560>

302 [Dainat, J., Hereñú, Da., Murray, K. D., Davis, E., Crouch, K., LucileSol, Agostinho, N.,](#)
303 [pascal-git, Zollman, Z. & tayyrov. \(2023\). NBISweden/AGAT: AGAT-v1.2.0 \(v1.2.0\).](#)
304 [Zenodo. <https://doi.org/10.5281/zenodo.8178877> Creators Jacques Dainat1 Darío](#)
305 [Hereñú Dr. K. D. Murray2 Ed Davis3 Kathryn Crouch4 LucileSol Nuno Agostinho5](#)
306 [pascal-git Zachary Zollman tayyrov Show affiliations 1. IRD 2. Max Planck Institut für](#)
307 [Entwicklungsbiologie, Tübingen 3. Oregon State University 4. @VEuPathDB 5.](#)
308 [European Bioinformatics Institute | EMBL-EBI. \(n.d.\). NBISweden/AGAT: AGAT-v1.2.0.](#)
309 <https://doi.org/10.5281/zenodo.8178877>

310 De Oliveira, A. L., Wollesen, T., Kristof, A., Scherholz, M., Redl, E., Todt, C., Bleidorn, C., &
311 Wanninger, A. (2016). Comparative transcriptomics enlarges the toolkit of known
312 developmental genes in mollusks. *BMC Genomics*, 17(1), 905.
313 <https://doi.org/10.1186/s12864-016-3080-9>

314 Elnaggar, A., Heinzinger, M., Dallago, C., Rehawi, G., Wang, Y., Jones, L., Gibbs, T., Feher,
315 T., Angerer, C., Steinegger, M., Bhowmik, D., & Rost, B. (2022). ProtTrans: Toward
316 Understanding the Language of Life Through Self-Supervised Learning. *IEEE*
317 *Transactions on Pattern Analysis and Machine Intelligence*, 44(10), 7112–7127.
318 <https://doi.org/10.1109/TPAMI.2021.3095381>

319 Erséus, C., Williams, B. W., Horn, K. M., Halanych, K. M., Santos, S. R., James, S. W.,
320 Creuzé des Châtelliers, M., & Anderson, F. E. (2020). Phylogenomic analyses reveal a
321 Palaeozoic radiation and support a freshwater origin for clitellate annelids. *Zoologica*

322 *Scripta*, 49(5), 614–640. <https://doi.org/10.1111/zsc.12426>

323 Fernández, R., & Gabaldón, T. (2020). Gene gain and loss across the metazoan tree of life.
324 *Nature Ecology & Evolution*, 4(4), 524–533. <https://doi.org/10.1038/s41559-019-1069-x>

325 Fernández, R., Tonzo, V., Simón Guerrero, C., Lozano-Fernandez, J., Martínez-Redondo,
326 G. I., Balart-García, P., Aristide, L., Eleftheriadi, K., & Vargas-Chávez, C. (2022).
327 MATEdb, a data repository of high-quality metazoan transcriptome assemblies to
328 accelerate phylogenomic studies. *Peer Community Journal*, 2(e58).
329 <https://doi.org/10.24072/pcjournal.177>

330 [Grabherr, M.G., Haas, B.J., Yassour, M., Levin, J.Z., Thompson, D.A., Amit, I., Adiconis, X.,
331 et al. 2011. Full-Length Transcriptome Assembly from RNA-Seq Data without a
332 Reference Genome. *Nature Biotechnology*. 29 \(7\), 644-52.
333 <https://doi.org/10.1038/nbt.1883>](#)

334 [Krug, P.J., Caplins, S.A., Algosó, K., Thomas, K., Valdés, A.A., Wrade, R., Wong, N.L.W.S.,
335 et al. \(2022\). Phylogenomic resolution of the root of Panpulmonata, a hyperdiverse
336 radiation of gastropods: new insights into the evolution of air breathing. *Proc. R. Soc. B*,
337 \[289, 20211855. <http://doi.org/10.1098/rspb.2021.1855>\]\(#\)](#)

338 Littmann, M., Heinzinger, M., Dallago, C., Olenyi, T., & Rost, B. (2021). Embeddings from
339 deep learning transfer GO annotations beyond homology. *Scientific Reports*, 11(1),
340 1160. <https://doi.org/10.1038/s41598-020-80786-0>

341 Manni, M., Berkeley, M. R., Seppey, M., Simão, F. A., & Zdobnov, E. M. (2021). BUSCO
342 Update: Novel and Streamlined Workflows along with Broader and Deeper Phylogenetic
343 Coverage for Scoring of Eukaryotic, Prokaryotic, and Viral Genomes. *Molecular Biology
344 and Evolution*, 38(10), 4647–4654. <https://doi.org/10.1093/molbev/msab199>

345 Martínez-Redondo, G.I., Barrios-Núñez, I., Vázquez-Valls, M., Rojas, A.M., &
346 Fernández, R. (2024). Illuminating the functional landscape of the dark proteome across
347 the Animal Tree of Life through natural language processing models.

348 Osmanski, A. B., Paulat, N. S., Korstian, J., Grimshaw, J. R., Halsey, M., Sullivan, K. A. M.,
349 Moreno-Santillán, D. D., Crookshanks, C., Roberts, J., Garcia, C., Johnson, M. G.,

350 Densmore, L. D., Stevens, R. D., Zoonomia Consortium†, Rosen, J., Storer, J. M.,
351 Hubley, R., Smit, A. F. A., Dávalos, L. M., ... Ray, D. A. (2023). Insights into mammalian
352 TE diversity through the curation of 248 genome assemblies. *Science*, 380(6643),
353 eabn1430. <https://doi.org/10.1126/science.abn1430>

354 Sayers, E. W., Bolton, E. E., Brister, J. R., Canese, K., Chan, J., Comeau, D. C., Farrell, C.
355 M., Feldgarden, M., Fine, A. M., Funk, K., Hatcher, E., Kannan, S., Kelly, C., Kim, S.,
356 Klimke, W., Landrum, M. J., Lathrop, S., Lu, Z., Madden, T. L., ... Sherry, S. T. (2023).
357 Database resources of the National Center for Biotechnology Information in 2023.
358 *Nucleic Acids Research*, 51(D1), D29–D38. <https://doi.org/10.1093/nar/gkac1032>

359 [Shen, W., Sipos, B., & Zhao, L. \(2024\). SeqKit2: A Swiss army knife for sequence and](#)
360 [alignment processing. *iMeta*, e191. <https://doi.org/10.1002/imt2.191>](#)

361 Song, H., Wang, Y., Shao, H., Li, Z., Hu, P., Yap-Chiongco, M. K., Shi, P., Zhang, T., Li, C.,
362 Wang, Y., Ma, P., Vinther, J., Wang, H., & Kocot, K. M. (2023). Scaphopoda is the sister
363 taxon to Bivalvia: Evidence of ancient incomplete lineage sorting. *Proceedings of the*
364 *National Academy of Sciences of the United States of America*, 120(40), e2302361120.
365 <https://doi.org/10.1073/pnas.2302361120>

366 Thoma, M., Missbach, C., Jordan, M. D., Grosse-Wilde, E., Newcomb, R. D., & Hansson, B.
367 S. (2019). Transcriptome surveys in silverfish suggest a multistep origin of the insect
368 odorant receptor gene family. *Frontiers in Ecology and Evolution*, 7.
369 <https://doi.org/10.3389/fevo.2019.00281>

370 UniProt Consortium. (2023). UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic*
371 *Acids Research*, 51(D1), D523–D531. <https://doi.org/10.1093/nar/gkac1052>

372 [Weisman, C.M., Murray, A.W., & Eddy, S.R. \(2022\). Mixing genome annotation](#)
373 [methods in a comparative analysis inflates the apparent number of lineage-specific](#)
374 [genes. *Current Biology*, 32 \(12\), 2632-2639.e2.](#)
375 <https://doi.org/10.1016/j.cub.2022.04.085>

376 Zapata, F., Wilson, N. G., Howison, M., Andrade, S. C. S., Jörger, K. M., Schrödl, M., Goetz,
377 F. E., Giribet, G., & Dunn, C. W. (2014). Phylogenomic analyses of deep gastropod

378 relationships reject Orthogastropoda. *Proceedings. Biological Sciences / The Royal*
379 *Society*, 281(1794), 20141739. <https://doi.org/10.1098/rspb.2014.1739>