

1 **FAIRification of DMRichR Pipeline: Advancing Epigenetic Research on Environmental and**
2 **Evolutionary Model Organisms**

3

4

5 Wassim Salam^{1,2} <https://orcid.org/0009-0001-0372-195X>

6 Marcin W. Wojewodzc^{2,3,*} <https://orcid.org/0000-0003-2501-5201>

7 Dagmar Frisch^{4,*} <https://orcid.org/0000-0001-9310-2230>

8

9

10

11 ¹Department of Mathematics and Computer Science, Freie Universität Berlin, Berlin,
12 Germany

13 ²Cancer Registry of Norway, Norwegian Institute of Public Health, Oslo, Norway

14 ³Department of Chemical Toxicology, Norwegian Institute of Public Health, Oslo, Norway

15 ⁴Department of Evolutionary and Integrative Ecology, Leibniz Institute of Freshwater
16 Ecology and Inland Fisheries, Berlin, Germany

17

18 *Corresponding authors:

19 dagmar.frisch@igb-berlin.de; MarcinWlodzimierz.Wojewodzc@fhi.no

20

21

22

23 **Key words:** Ecology, Epigenetics, *Daphnia pulex*, DMRichR, Methylation, Non-model
24 organism, risk assessment

25

26 Abstract

27 Bioinformatics tools often prioritize humans or human-related model organisms,
28 overlooking the requirements of environmentally relevant species, which limits their use in
29 ecological research. This gap is particularly challenging when implementing existing
30 software, as inadequate documentation can delay the innovative use of environmental
31 models for modern risk assessment of chemicals that can cause aberration in methylation
32 patterns. The establishment of fairness in ecological and evolutionary studies is already
33 constrained by more limited resources in these fields of study, and an additional imbalance
34 in tool availability further hinders comprehensive ecological research. To address these
35 gaps, we adapted the DMRichR package, a tool for epigenetic analysis, for use with custom,
36 non-model genomes. As an example we here use the crustacean *Daphnia*, a keystone grazer
37 in aquatic ecosystems. This adaptation involved the modification of specific code,
38 computing three new species-specific packages (BSgenome, TxDb, and org.db), and
39 computing a CpG islands track using the makeCGI package. Additional adjustments to the
40 DMRichR package were also necessary to ensure proper functionality. The developed
41 workflow can now be applied not only to different *Daphnia* species that were previously
42 unsupported, but also to any other species for which an annotated reference genome is
43 available.

44

45 **Introduction**

46 Epigenetic research is a crucial field in ecological research for understanding how organisms
47 adapt to environmental challenges (Thiebaut, Hemerly and Ferreira 2019; McGuigan,
48 Hoffmann and Sgrò 2021; Lamka *et al.* 2022) as well as for application in ecotoxicological
49 studies that involve non-model species (Vandegheuchte and Janssen 2014; Šrut 2021).
50 However, the related bioinformatics tools are predominantly oriented towards humans or
51 model organisms for humans, often neglecting the requirements for environmental model
52 organisms. This disparity hinders their development and application in ecological research
53 or modern risk assessment for chemicals. Additionally, a lack of clear documentation or the
54 absence of necessary dependencies further complicates the implementation of existing
55 software for these organisms.

56 Differential genome-wide methylation analysis involves the comparison of methylation
57 patterns across different conditions to understand the impacts on gene regulation and
58 consequently gene expression (Parle-Mcdermott and Harrison 2011; Li and Tollefsbol 2021).
59 Currently, a multitude of tools exists that are suitable for Differentially Methylated Region
60 (DMR) analysis, each of which makes use of a different differential methylation test. These
61 include but are not limited to: Fisher's exact test, BSmooth, MethylKit, MethylSig, DSS,
62 Metilene, RADMeth, Biseq. Each tool uses a unique approach, with none of them
63 consistently outperforming others during benchmark testing (Piao *et al.* 2021).

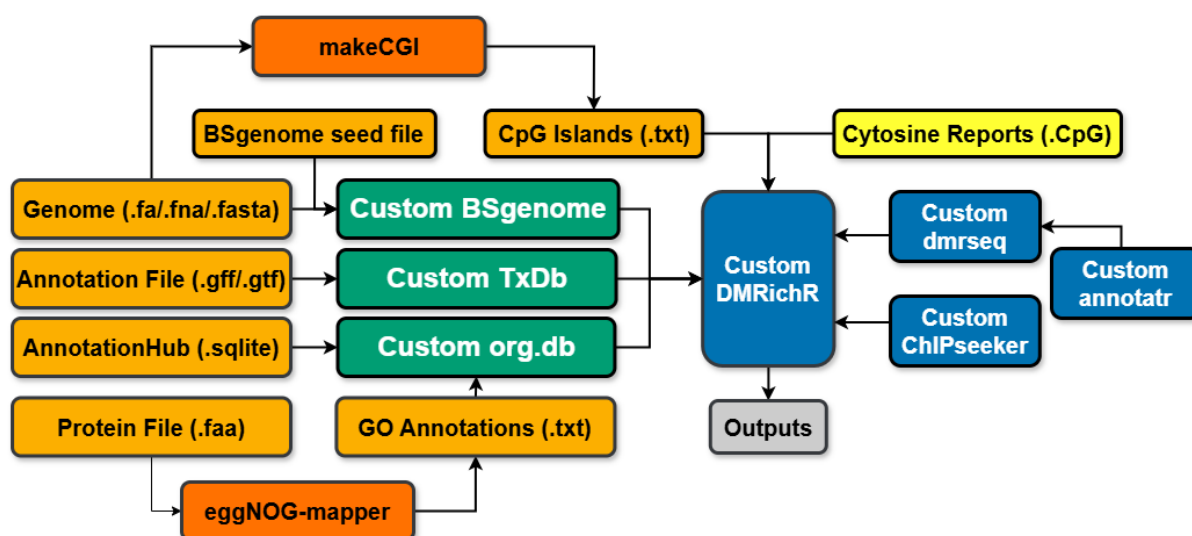
64 One recent tool that stands out for its robust use of a Bayesian framework is DMRichR
65 (v1.7.8) (Hansen, Langmead and Irizarry 2012; Korthauer *et al.* 2019; Laufer *et al.* 2021). It is
66 a powerful resource in epigenetic research which contrasts with other DMR analysis tools by
67 combining both the dmrseq (v1.15.1) and bsseq (v1.38.0) algorithms to identify
68 differentially methylated regions (DMRs) with greater precision and accuracy. It improves
69 analysis by incorporating prior knowledge and probabilistic models to better capture true
70 methylation changes, thus reducing false positives. This method prioritizes methylation sites
71 with higher coverage, using them to infer nearby sites with less data. To identify DMRs, this
72 approach focuses on the comparison of groups of samples (e.g. treatments) rather than
73 within-group levels, and allows the analysis of samples even at a lower sequencing depth of
74 1-5x (Laufer 2023). DMRs are identified in two steps: first, pooling and weighting data from
75 high-coverage sites to detect differences, and then statistically testing these regions to
76 identify significant changes across the genome (Korthauer *et al.* 2019; Laufer 2023). The

77 biggest notable drawback of DMRichR is that it is primarily available for typical model
78 species and its application can be challenging for researchers with limited bioinformatics
79 experience. As of its last update in October 2023, only 15 species are supported, which
80 include *Homo sapiens* and *Mus musculus*, while excluding many environmentally important
81 species or evolutionary models such as *Daphnia*.

82 One of such important additions to the suite of taxa to be used with DMRichR are aquatic
83 organisms. *Daphnia* are keystone organisms in aquatic ecosystems, serving as a vital link in
84 the food web between primary producers and higher trophic levels (Miner *et al.* 2012;
85 Ogorelec *et al.* 2020). In the last decade, the planktonic crustacean *Daphnia* has emerged as
86 an important model in ecological and evolutionary research, and is supported by institutions
87 such as the National Institutes of Health (NIH) (Ebert 2005, 2022). Despite their extensive
88 use in genetic, ecotoxicological, and ecological research (Colbourne *et al.* 2022), genomic
89 resources are typically built for *D. pulex* and *D. magna*, but even for these, bioinformatic
90 tools remain scarce. Providing the implementation of DMRichR for these taxa will facilitate
91 access to the pipelines in the research community, and could potentially drive modern
92 chemical risk assessment when epigenetic effects are of interest.

93 The use of custom genomes with the R package DMRichR can be achieved by making
94 appropriate modifications in the code of annotationDatabases.R (Laufer 2023). This requires
95 in a first step the computation of three new packages (Fig.1): BSgenome, TxDb and org.db
96 (for description of functionality see below). Next, a CpG islands track must be computed,
97 which can be done using the makeCGI package (Irizarry, Wu and Feinberg 2009; Wu *et al.*
98 2010). Lastly, a few modifications have to be made to packages used within DMRichR's
99 code, namely within Dmrseq, Annotatr and CHIPseeker.

100



101

102

103 **Figure 1.** Steps required to adapt DMRichR to a new organism with an annotated genome.104 **Yellow:** User input files (methylation calls) resulting from a Whole-Genome Bisulfite105 Sequencing (WGBS) experiment; **Light-orange:** User input files with information on the106 reference genome; **Dark-orange:** Intermediary tools for producing input files; **Green:** Newly-107 computed species-specific packages; **Blue:** Modified versions of DMRichR and additional108 packages used within it; **Grey:** Results generated by DMRichR, which include Blocks,

109 Differentially Methylated Regions (DMRs), Smoothed Individual Methylation Values,

110 Heatmaps.

111

112

113 **Implementation**

114 An overview of the workflow is provided in Fig. 1. While here we executed the workflow

115 specifically for *Daphnia pulex*, it could also be applied more generally for any other species

116 with an annotated reference genome following the same steps. The R code (v4.3.2) (R Core

117 Team 2023) associated with each package's computation will be published in Appendix 1. In

118 the next sections we describe the generation or modification made to provide all packages

119 and steps for DMRichR to function with the *Daphnia pulex* genome and the cytosine report

120 files produced by Bismark (Krueger and Andrews 2011).

121

122 *BSgenome*

123 The BSgenome package enables users to efficiently manage and analyze whole genome

124 sequences. It provides tools for tasks such as extracting genomic sequences and conducting
125 genome-wide analyses. Apart from DNA methylation analysis, this has additional value for
126 use in other applications such as sequence alignment, variant calling, histone modification
127 (ChIPseq) analysis, and motif discovery (Pagès 2024). The first step is to write a "seed file"
128 following BSgenome guidelines (Pagès 2020), which contains package metadata and
129 instructions on how the package should be compiled. The *D. pulex* ASM2113471v1 genome
130 assembly (NCBI 2021) was used to produce the seed file
131 BSgenome.Dpulex.NCBI.ASM2113471v1-seed (Appendix 2). Following recommended
132 naming conventions, the package BSgenome.Dpulex.NCBI.ASM2113471v1 was computed
133 (Appendix 3).

134

135 *TxDb*

136 Using either a GFF (General Feature Format) or GTF (Gene Transfer Format) file, a transcript
137 annotation package can be computed. It contains transcript-related annotations (like exons,
138 introns, UTRs), which in addition to being a central package for DMRichR, may be useful for
139 applications such as transcriptomics, RNA-Seq data analysis, Chip-Seq annotation, and gene
140 structure analysis. The *D. pulex* GTF file (NCBI 2022) was used, and in a two-step process,
141 the package TxDb.Dpulex.NCBI.ASM2113471v1.knownGene was computed (Appendix 4).

142

143 *org.db*

144 This package contains mappings between a central identifier (e.g., Entrez gene IDs) and
145 other identifiers (e.g. gene symbol, gene name, gene ontology, chromosome)(Morgan and
146 Arora 2014). The AnnotationHub package already contains a *D. pulex* database in SQLite
147 format, with "GID" (Entrez ID) as a central key for its tables. A GO (Gene Ontology)
148 annotation file was obtained by passing the *D. pulex* protein file (protein.faa.gz)(NCBI 2022)
149 to eggNOG-mapper (Huerta-Cepas et al. 2019; Cantalapiedra et al. 2021) with default
150 options. By combining these components we computed org.Duplex.eg.db, following org.db
151 package naming conventions (Appendix 5).

152

153 *CpG Islands*

154 By using the makeCGI package (Irizarry, Wu and Feinberg 2009; Wu et al. 2010), CpG Islands
155 of any available annotated genome can be de novo discovered to build a CpG islands track,

156 which is an integral component for the functionality of the DMRichR analysis. Computation
157 of the BSgenome package is a prerequisite to this step (Fig. 1), as makeCGI loads the
158 specified genome from the latter. The posterior probability is one of many parameters that
159 the user can modify, which affects how the package decides what defines a CpG Island.
160 From the CGI files of different organisms already available on the Hao Wu Lab website
161 (Appendix 6), a posterior probability of 0.99 was chosen for all genomes except for that of
162 the fruit fly *Drosophila melanogaster* (0.975). makeCGI was executed with
163 BSgenome.Dpulex.NCBI.ASM2113471v1 (Appendix 7), with a chosen posterior probability
164 similar to that of *D. melanogaster*, because the *Daphnia* genome has a higher resemblance
165 to that of the fruit fly than to the other listed genomes. The genome of *Daphnia*, like that of
166 the fruit fly, is characterized by a small number of methylated bases (Asselman et al. 2016;
167 Kusari et al. 2017). A text file containing CpG Islands entries was therefore produced
168 (Appendix 8). These parameters might need further tuning for other species.

169

170 **Modifications of DMRichR and Additional Packages**

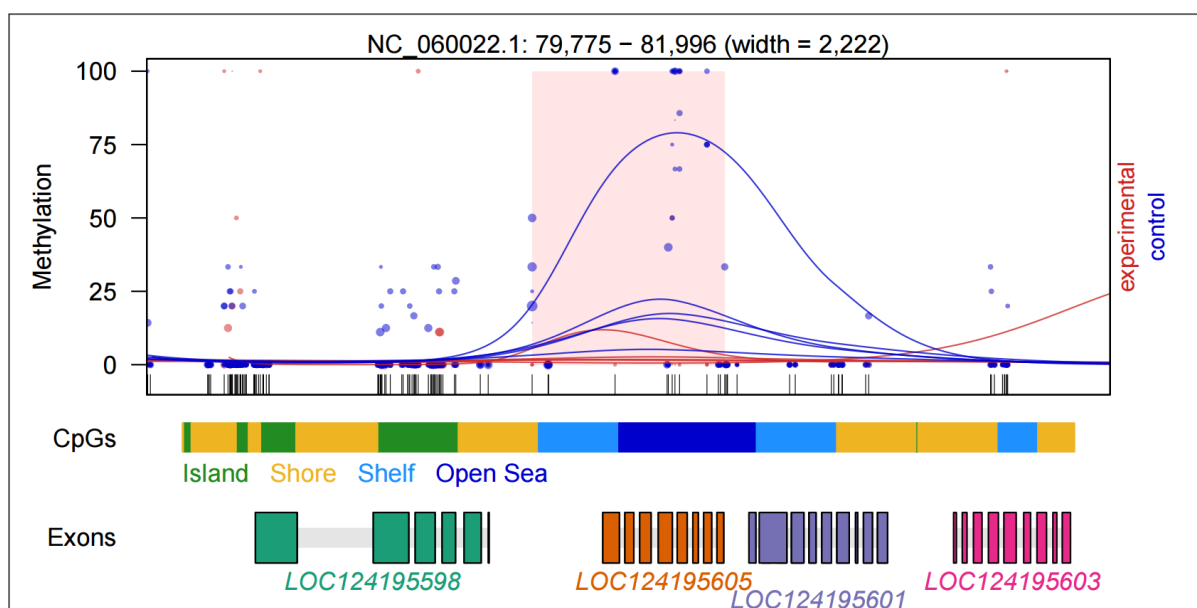
171 Multiple snippets of code were adjusted within the DMRichR package. Below we provide a
172 short overview of the changes, which can be viewed in detail in Appendix 9 :

- 173 • Added new genome “Dpulex” and integrated CGI annotations
- 174 • Integrated new BSgenome, TxDb and org.db packages in “annotationDatabases.R”
- 175 • Modifications were made to arguments such as “minInSpan”, “bpSpan”,
176 “maxGapSmooth”, “maxGap”, “minNumRegion” and “blockSize”. This was done because
177 *D. pulex* has not only a significantly smaller size but also low and sparse methylation.
178 Lowering the threshold of default arguments allows for more methylation data to be
179 captured.
- 180 • Added argument “cytosineReportFormat” (default NULL). Setting a value of “nf-
181 core/methylseq” (Ewels et al. 2023) would enable DMRichR to process cytosine reports
182 generated by nf-core/methylseq which are produced with a slightly different naming
183 convention than when produced by Bismark.
- 184 • The packages dmrseq, annotatr and ChIPseeker were adjusted to allow the integration
185 of the new *D. pulex* genome. The respective changes can be seen in Appendix 10, 11 and
186 12.

187

188 **Case Study**

189 Proper functionality of the modified R package DMRichR is demonstrated using sample data
 190 from an unpublished study involving Whole-Genome Bisulfite Sequencing data from
 191 *Daphnia pulicaria*, a member of the *Daphnia pulex* species complex (Dudycha and Tessier
 192 1999). The provided input example presents a cytosine report generated from 10 samples (5
 193 control vs. 5 experimental), and for demonstration purposes has been reduced to report
 194 methylation only on chromosome NC_060022.1. The code for this test run is provided in
 195 Section 4, which contains instructions on setup and installation. Detailed instructions about
 196 how to test the customized DMRichR package can be found in **Appendix 13**. The cytosine
 197 reports of each sample (Appendix 14), which are used as input for DMRichR, were produced
 198 by nf-core/methylseq (v2.4.0) (Ewels *et al.* 2023). The DMR plot shown in Fig. 2 below
 199 displays one of many DMRs obtained by this case study run, which successfully
 200 demonstrates the FAIRification of the DMRichR pipeline.



201

202 **Figure 2.** DMR plot displaying a DMR consisting of 13 CpGs with 23% hypomethylation in
 203 experimental samples compared to control samples. The methylation level of each CpG site
 204 in an individual sample is shown as a point, with its size directly proportional to its coverage.
 205 Smoothed methylation levels are represented by lines, color-coded as blue for control
 206 samples and red for experimental samples. A track of CpG and gene annotations are
 207 additionally displayed under the plot, retrieved from the computed CGI track and org.db
 208 package respectively.

209 **Data and Code availability**

210 All analyses were performed using R Statistical Software (v4.3.2) (R Core Team 2023) and
211 Bioconductor (v3.18) (Huber *et al.* 2015). The BSgenome seed-file for *D. pulex*, the code
212 used to compute the above-mentioned packages and the CpG islands list, as well as the
213 sample data (cytosine reports) used in the DMRichR test-run will be made publicly available
214 upon publication. For ease of use, the installations of the computed *D. pulex* packages
215 (BSgenome, TxDb and org.db) were seamlessly integrated into the custom DMRichR
216 package. However, they are as well readily available to use independently of DMRichR. This
217 can prove useful for specific applications, some of which are mentioned in the respective
218 package sections.

219

220 **Conclusion**

221 Integrating support for the *D. pulex* genome into the DMRichR package represents a
222 significant advancement in the field of ecological and evolutionary genomics. It strengthens
223 the capacity for high-resolution analysis of DNA methylation patterns in *D. pulex*, and
224 enhances the possibility of using whole genome methylation in a modern risk assessment
225 for chemicals. The incorporation of support for the *D. pulex* genome to DMRichR thus allows
226 researchers to leverage this tool's robust functionalities to investigate epigenetic
227 modifications and efficiently use sparse information across different treatments with
228 greater precision. This adaptation not only facilitates deeper insights into the adaptive
229 mechanisms and environmental responses of *D. pulex* but also creates possible use for risk
230 assessment using epigenetics that is still underexplored in ecological studies.

231 The workflow described here sets a precedent for similar enhancements in other species.

232 The process involves the careful annotation of the target species' genome, followed by
233 integration into the DMRichR framework, thereby enabling the broader scientific
234 community to extend these powerful analytical capabilities to a diverse array of organisms.
235 With an annotated reference genome, increasingly available for many non-model species,
236 the workflow we have described and tested here is particularly beneficial for researchers
237 studying methylation patterns in ecologically and evolutionary significant species, as it
238 bridges the gap between advanced bioinformatics tools and ecological research, fostering a
239 more comprehensive understanding of epigenetic regulation in varied environmental
240 contexts. Lastly, the packages produced in this work contribute not only to the

241 advancement of differential methylation analysis, but also to other applications improving
242 FAIRness of these tools for environmental research.

243

244 **Acknowledgments**

245 For computing time we would like to thank the High-Performance Computing Service of
246 ZEDAT, Freie Universität Berlin.

247

248 **Funding**

249 DF acknowledges funding by the Deutsche Forschungsgemeinschaft (DFG, German Research
250 Foundation – Project number 461099895) and from the European Union's Horizon 2020
251 research and innovation program under the Marie Skłodowska-Curie grant agreement No.
252 658714. MWW acknowledges funding from Marie Skłodowska Curie action FP7-PEOPLE-
253 2013-IEF (EU, GB). DF and MWW acknowledge NERC Biomolecular Analysis Facility
254 (Liverpool, GB, NBAF998). The bioinformatic analysis was partially funded by the German
255 Federal Ministry of Education and Research (BMBF, Förderkennzeichen 033W034A).

256

257 **Competing Interests**

258 None declared.

259

260 **Availability and Implementation:** Code and data will be made available after peer review,
261 upon publication in github and Zenodo.

262

263 **Author Contributions**

264 WS: Data Curation - lead, Investigation - lead, Writing - original draft; DF & MW:
265 Conceptualisation - equal, Investigation - equal, Supervision - equal, Writing; review and
266 editing - equal; Data Production - equal

267

268 **References**

- 269 Colbourne JK, Shaw JR, Sostare E et al. Toxicity by descent: A comparative approach for
270 chemical hazard assessment. *Environ Adv* 2022;9:100287.
- 271 Dudycha JL, Tessier AJ. Natural genetic variation of life span, reproduction and juvenile
272 growth in *Daphnia*. *Evolution* 1999;53:1744–56.
- 273 Ebert D. Introduction to *Daphnia* Biology. Ecology, Epidemiology, and Evolution of
274 Parasitism in *Daphnia* [Internet]. National Center for Biotechnology Information (US),
275 2005.
- 276 Ebert D. *Daphnia* as a versatile model system in ecology and evolution. *EvoDevo* 2022;13:16.
- 277 Ewels P, Hüther P, Sateesh P et al. nf-core/methylseq: [2.4.0] Gillespie Gaia. 2023, DOI:
278 10.5281/zenodo.8029942.
- 279 Hansen KD, Langmead B, Irizarry RA. BSmooth: from whole genome bisulfite sequencing
280 reads to differentially methylated regions. *Genome Biol* 2012;13:R83.
- 281 Huber W, Carey VJ, Gentleman R et al. Orchestrating high-throughput genomic analysis with
282 Bioconductor. *Nat Methods* 2015;12:115–21.
- 283 Irizarry RA, Wu H, Feinberg AP. A species-generalized probabilistic model-based definition of
284 CpG islands. *Mamm Genome Off J Int Mamm Genome Soc* 2009;20:674–80.
- 285 Korthauer K, Chakraborty S, Benjamini Y et al. Detection and accurate false discovery rate
286 control of differentially methylated regions from whole genome bisulfite sequencing.
287 *Biostat Oxf Engl* 2019;20:367–83.
- 288 Krueger F, Andrews SR. Bismark: a flexible aligner and methylation caller for Bisulfite-Seq
289 applications. *Bioinformatics* 2011;27:1571–2.
- 290 Lamka GF, Harder AM, Sundaram M et al. Epigenetics in Ecology, Evolution, and
291 Conservation. *Front Ecol Evol* 2022;10.
- 292 Laufer B. Introduction to DMRichR. 2023.
293 <https://www.benlaufer.com/DMRichR/articles/DMRichR.html>
- 294 Laufer B, Hwang H, Jianu JM et al. Low-pass whole genome bisulfite sequencing of neonatal
295 dried blood spots identifies a role for RUNX1 in Down syndrome DNA methylation
296 profiles. *Hum Mol Genet* 2021;29:3465–76.
- 297 Li S, Tollefsbol TO. DNA methylation methods: Global DNA methylation and methylomic
298 analyses. *Methods San Diego Calif* 2021;187:28–43.
- 299 McGuigan K, Hoffmann AA, Sgrò CM. How is epigenetics predicted to contribute to climate

- 300 change adaptation? What evidence do we need? *Philos Trans R Soc Lond B Biol Sci*
301 2021;376:20200119.
- 302 Miner BE, De Meester L, Pfrender ME et al. Linking genes to communities and ecosystems:
303 *Daphnia* as an ecogenomic model. *Proc R Soc B Biol Sci* 2012;279:1873–82.
- 304 Ogorelec Z, Wunsch C, Kunzmann A et al. Large daphniids are keystone species that link fish
305 predation and phytoplankton in trophic cascades. *Fundam Appl Limnol Arch Für*
306 *Hydrobiol* 2020;194, DOI: 10.1127/fal/2020/1344.
- 307 Pagès H. BSgenome: Software infrastructure for efficient representation of full genomes and
308 their SNPs. *Bioconductor* 2024.
- 309 Parle-Mcdermott A, Harrison A. DNA Methylation: A Timeline of Methods and Applications.
310 *Front Genet* 2011;2, DOI: 10.3389/fgene.2011.00074.
- 311 Piao Y, Xu W, Park KH et al. Comprehensive Evaluation of Differential Methylation Analysis
312 Methods for Bisulfite Sequencing Data. *Int J Environ Res Public Health* 2021;18:7975.
- 313 R Core Team. R: A language and environment for statistical computing. 2023.
- 314 Šrut M. Ecotoxicological epigenetics in invertebrates: Emerging tool for the evaluation of
315 present and past pollution burden. *Chemosphere* 2021;282:131026.
- 316 Thiebaut F, Hemerly AS, Ferreira PCG. A Role for Epigenetic Regulation in the Adaptation
317 and Stress Responses of Non-model Plants. *Front Plant Sci* 2019;10, DOI:
318 10.3389/fpls.2019.00246.
- 319 Vandegehuchte MB, Janssen CR. Epigenetics in an ecotoxicological context. *Mutat Res*
320 *Genet Toxicol Environ Mutagen* 2014;764–765:36–45.
- 321 Wu H, Caffo B, Jaffee HA et al. Redefining CpG islands using hidden Markov models.
322 *Biostatistics* 2010;11:499–514.