



26 **Data availability**

27 The ‘pynnotate’ public repository is available at <https://github.com/fernandacaron/pynnotate>.

28

29 **Conflict of interest**

30 The authors have no conflicts of interest to declare.

31

32 **Author contribution statements**

33 FSC (Conceptualisation, Methodology, Software, Validation, Formal analysis, Investigation,

34 Writing—original draft, Writing—review & editing), FMM (Conceptualisation, Methodology,

35 Software, Validation, Formal analysis, Investigation, Writing—review & editing), MMAS

36 (Conceptualisation, Methodology, Validation, Writing—original draft, Writing—review &

37 editing), and FMCBD (Conceptualisation, Methodology, Validation, Writing—review &

38 editing, Project Administration).

39

40 **Abstract**

41 Pynnotate is a Python-based tool designed for automated retrieval, parsing, and extraction of  
42 annotated gene sequences from GenBank records. The tool addresses the common challenges  
43 researchers face when working with GenBank data, including inconsistent gene nomenclature,  
44 redundant sequences, and the need for standardised gene extraction across multiple taxa.  
45 Pynnotate operates through both a graphical user interface and a command-line interface,  
46 making it accessible to users with varying levels of bioinformatics experience. The tool  
47 supports flexible sequence retrieval through manually defined accession numbers or NCBI  
48 query terms, and offers three distinct filtering modes: unconstrained (all sequences), strict (one  
49 sequence per species prioritising gene completeness), and flexible (multiple sequences per  
50 species when contributing different genes). Key features include synonym resolution for gene  
51 names, customizable sequence headers, metadata tracking, and automated gene extraction into  
52 separate files. Built-in dictionaries support animal and plant mitochondrial DNA, chloroplast  
53 DNA, and ribosomal DNA, and allow users to provide custom synonym dictionaries. The tool  
54 generates structured output including FASTA files, metadata matrices, and detailed logs,  
55 facilitating integration with downstream analyses. Designed for speed and scalability,  
56 pynnotate efficiently handles large datasets, allowing quick retrieval and extraction of  
57 annotated sequences across multiple taxa. Finally, pynnotate serves as a valuable resource for  
58 both research applications and educational settings, particularly benefiting educators  
59 conducting bioinformatics analyses with students with limited command-line experience.

60 **Keywords:** bioinformatics; comparative genomics; feature extraction; molecular evolution;  
61 phylogenetics; Python; sequence annotation.

62

## 63 1. Introduction

64 Public molecular databases have revolutionised molecular evolution research by  
65 providing access to vast amounts of genomic data across diverse taxa. Pioneering the field,  
66 Margaret Dayhoff and her colleagues were the first to publish an online database of molecular  
67 sequences in 1965, aiming to assist research and make protein data easily accessible (Hagen,  
68 2011). This *Atlas of Protein Sequence and Structure* was particularly important because it  
69 inspired the creation of GenBank. Since its first release in 1982, the amount of sequence data  
70 in GenBank has roughly doubled every 18 months, reflecting the increasing scale of molecular  
71 data generation worldwide (<https://www.ncbi.nlm.nih.gov/genbank/statistics/>). Currently,  
72 GenBank hosts hundreds of millions of publicly available nucleotide sequences, but despite  
73 this wealth of data, retrieving and organising specific genetic features for targeted evolutionary  
74 analyses remains challenging, particularly for non-model organisms.

75 A key contribution of GenBank to molecular research is the availability of stable  
76 accession numbers and standardised records, which facilitate data traceability and  
77 reproducibility. Making molecular data publicly available has enabled collaborative science  
78 and has allowed researchers to replicate previously published results (Leray et al., 2019).  
79 However, to successfully replicate molecular evolution research, the processing and  
80 manipulation of sequences must be performed almost exactly to achieve a comparable  
81 alignment (Baykal et al., 2024). An issue arises when large numbers of these sequences must  
82 be analysed simultaneously because GenBank records are not necessarily standardised. For  
83 instance, different sequences may use various names for the same genes (e.g., the cytochrome  
84 oxidase I may be named COI, or COX1, or CO1), or sequences related to one feature might be  
85 mistakenly registered under another feature's name. While it is standard to address these  
86 problems when curating datasets for molecular evolution research projects, the specific  
87 methods used to handle these issues are not always clearly defined, limiting replicability.

88           There is a growing demand for flexible tools that streamline the retrieval of sequence  
89 features for evolutionary analysis. It is common practice in molecular evolutionary studies to  
90 use complete mitogenomes or similar records with annotated features, analysing these features  
91 separately. Currently, GenBank does not provide an option to download each feature from these  
92 records separately, other than manually. This limitation becomes a major bottleneck in large-  
93 scale comparative studies involving hundreds or thousands of annotated sequences. Moreover,  
94 the tools available for separate downloads by feature have several limitations that may prevent  
95 the full utilisation of GenBank data, such as not retrieving reference sequences (e.g., Borstein  
96 & O’Meara, 2018), potentially resulting in incomplete datasets.

97           Beyond research applications, there is growing interest in using bioinformatics tools for  
98 educational purposes, particularly those that combine graphical and command-line interfaces  
99 to support progressive instruction across different levels of expertise. Such tools have proven  
100 effective in promoting active learning and improving bioinformatics and evolutionary biology  
101 literacy among students with limited programming experience (Emery & Morgan, 2017; Harris  
102 et al., 2020). Instructors can employ these platforms to teach essential skills—such as accessing  
103 public databases, resolving gene-name inconsistencies, curating metadata, and extracting  
104 annotated features—while introducing broader conceptual topics in molecular evolution,  
105 including homology, alignment quality, and annotation errors. The integration of reproducible  
106 workflows into teaching environments also provides a valuable opportunity to emphasise open  
107 science practices and data transparency (e.g., Cokelaer et al., 2023; Toelch & Ostwald, 2018).  
108 Consequently, user-friendly and reproducible bioinformatics tools can serve as powerful  
109 educational resources, bridging the gap between theoretical concepts and hands-on data  
110 analysis in molecular evolution.

111           To address these research and educational gaps, we developed pynnotate, a flexible  
112 Python-based tool that standardises the extraction and processing steps, allowing for

113 reproducible outputs across diverse datasets. Pynnotate is not intended to perform genome  
114 assembly or large-scale annotation pipelines. Instead, it facilitates targeted sequence retrieval  
115 and curation of already annotated loci, particularly mitochondrial (mtDNA), plastid (cpDNA),  
116 ribosomal (rDNA), and selected nuclear markers commonly used in molecular evolution and  
117 systematics. In addition to its command-line interface, pynnotate also includes a graphical  
118 version, enabling use by researchers and students without experience in Python, including in  
119 classroom settings. Inspired by similar tools and packages developed in the fields of ecology  
120 and evolution, we named our module pynnotate, a reference to the Python programming  
121 language and its function related to gene feature annotation. It supports batch processing of  
122 large datasets by retrieving and parsing multiple GenBank records sequentially, allowing users  
123 to extract multiple features across many species in a single run. Here, we describe the  
124 implementation of pynnotate and demonstrate its utility through easy-to-follow tutorials.  
125 Pynnotate is publicly available and can be easily adapted for various studies using GenBank  
126 data.

127

## 128 **2. Package overview**

### 129 *a. Overview*

130 Pynnotate is a Python-based tool designed for automated retrieval, parsing, and  
131 extraction of annotated gene sequences from GenBank records, with a focus on mitochondrial  
132 and chloroplast genomes. The module is designed for workflows requiring standardised gene  
133 extraction across multiple taxa, with support for synonym resolution, customizable headers,  
134 and metadata tracking. Pynnotate operates on both manually defined accession numbers and  
135 NCBI query terms, enabling flexible sequence retrieval. Its extraction method recognises and  
136 parses multiple annotation fields (e.g., gene, product, note), detecting aliases using synonym  
137 dictionaries, and organises the output by gene. Currently implemented on pynnotate are

138 dictionaries for animal and plant mtDNA, cpDNA, rDNA, and some animal nuclear genes.  
139 Although not complete, an advantage of pynnotate is its flexibility, which allows the user to  
140 provide a custom synonym dictionary, and its functionalities will apply to any other sequence.  
141 Specific logic is implemented for tRNA-Leu and tRNA-Ser genes to differentiate positional  
142 isoforms based on genomic context. Optional filters allow control over gene overlap, record  
143 length, and the number of individuals per species, making the tool suitable for low- and high-  
144 throughput genomic data processing. In this context, pynnotate has three current filtering  
145 modes: (1) download all sequences matching the search query (unconstrained mode); (2)  
146 download one individual per species, choosing those that maximize the number of genes  
147 downloaded per species (strict); and (3) downloads many individuals per species, but choosing  
148 only one individual per gene (i.e., supermatrix) (flexible). Limitations include dependency on  
149 the accuracy and completeness of GenBank annotations.

150

#### 151 *b. Implementation*

152 Pynnotate is written in Python ( $\geq 3.8$ ) and can be executed on Windows, macOS, or  
153 Linux. Pynnotate relies primarily on the Biopython library for parsing GenBank records, and  
154 several other external dependencies (<https://github.com/fernandacaron/pynnotate>). These  
155 dependencies are automatically installed and require no additional compilation steps. Bugs and  
156 technical issues can be reported directly through the GitHub repository's issues page.

157 In the command-line, all execution settings of pynnotate are specified via a YAML  
158 configuration file. This file defines the data source (accession numbers or search terms), desired  
159 genes, sequence filters (e.g., minimum length, unique species), and output options such as  
160 header customisation and log generation. All results are saved in a structured directory  
161 containing FASTA files, logs, and metadata matrices, facilitating integration with downstream  
162 analytical tools. In addition to storing extracted sequences, pynnotate logs all retrieved records

163 and explicitly reports those that were not extracted, together with the reason for exclusion (e.g.,  
164 missing annotations or applied filters). This ensures full traceability of dataset assembly and  
165 allows users to review which records contributed to the final output. As we recognise that  
166 running Python scripts can be a barrier for many users in bioinformatics, we have developed a  
167 graphical user interface (GUI) alongside the module, which users can download and install  
168 locally (on Windows, macOS, or Linux). This application allows all module functionalities to  
169 be executed through an intuitive and user-friendly interface. In this way, we hope this tool will  
170 serve as a flexible and accessible resource for researchers to retrieve data from GenBank in an  
171 automated and efficient manner, while also supporting the reproducibility of the results.  
172 Furthermore, we believe that the graphical interface is particularly valuable for educators who  
173 wish to conduct bioinformatics analyses in classroom settings with students who have no prior  
174 experience with command-line tools, coding, or bioinformatics.

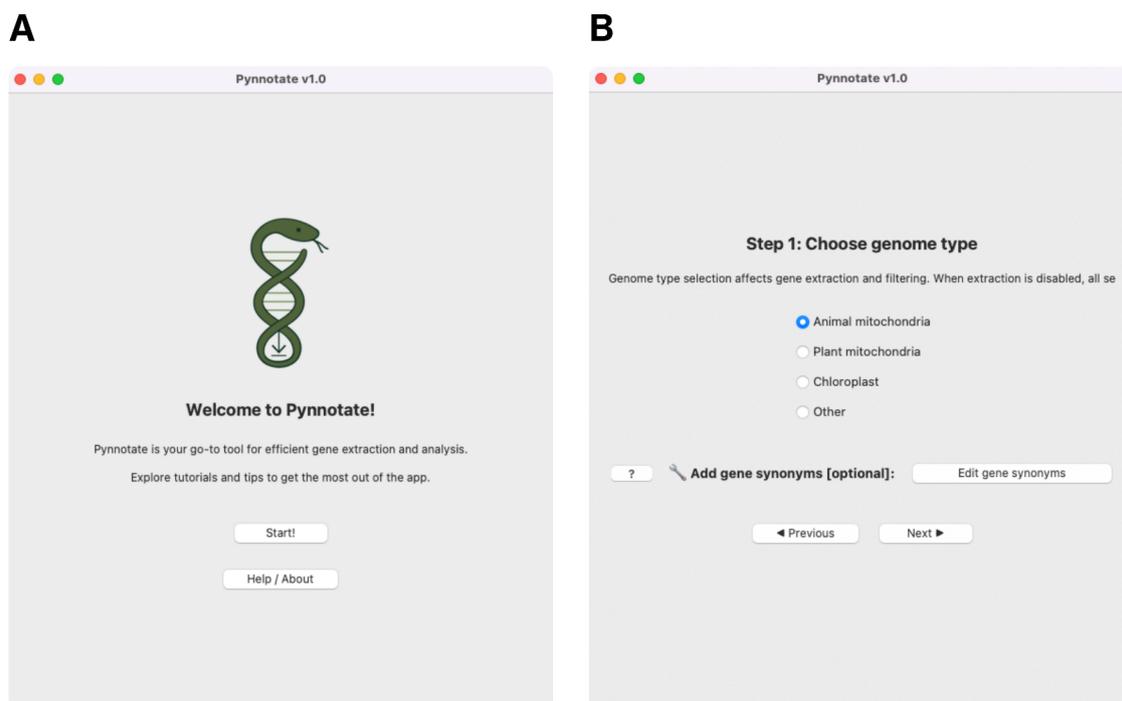
175

### 176 **3. Graphic interface workflow**

#### 177 *a. Getting started*

178 The pynnotate interface was designed to make it easy for users to get all functionalities  
179 in a simple step-by-step workflow (Figure 1A). The first required parameter to retrieve  
180 GenBank sequences is the genome type (Figure 1B). This genome type is used subsequently  
181 for gene extraction and filtering, so not all functionalities will be affected by this choice. More  
182 specifically, when gene extraction is disabled, all sequences matching your search will be  
183 downloaded regardless of genome type. The genome type can be one of (1) animal  
184 mitochondria, (2) plant mitochondria, (3) chloroplast, or (4) other. For the first three options,  
185 some synonyms are already built into pynnotate. These synonyms are used to identify which  
186 sequences in GenBank databases correspond to each gene. Since GenBank lacks gene names  
187 standardisation, not all gene synonyms are included in our dictionary. Therefore, users can

188 specify new synonyms using the "Add gene synonyms" option in Step 1 (Figure 1B). We  
189 recommend running the program first to identify any unrecognised gene synonyms.  
190 Alternatively, if options 1-3 do not meet the user's needs, a new dictionary can also be provided  
191 using the genome type "Other" and specifying which synonyms should be considered in the  
192 "Add gene synonyms" option.  
193



194  
195 **Figure 1.** Pynnotate graphic interface. **A.** Welcome page including citation information. **B.**  
196 Step 1: Defining genome type and gene synonyms.

197  
198 *b. Assembling a search query*

199 Pynnotate searches for and downloads sequences from GenBank either by assembling  
200 a search query with user-defined parameters or directly using accession numbers provided by  
201 the user. In Step 2 (Figure 2), users can specify these accession numbers in the first field.  
202 Alternatively, a search query can be assembled using the remaining available fields. When

203 choosing this approach, users can specify particular genes, organism names (e.g., families,  
204 genera, species), and publication terms. An additional field allows users to specify any other  
205 query components. Furthermore, users can restrict the query to mitochondrial genes, the  
206 mitogenome, or the chloroplast genome. As a side note, since GenBank does not explicitly  
207 support such restricted filtering, some sequences outside the intended scope may still be  
208 downloaded when using these options. Users can also choose to exclude unannotated sequences  
209 ('UNVERIFIED' in GenBank) at this step. Finally, all fields are combined into a single query.  
210 The search shown in Figure 2, for example, would generate '(COI) AND (Anura[Organism]).'  
211

The screenshot shows a web browser window titled "Pynnotate v1.0". The main heading is "Step 2: Build custom search query". Below this, there are several input fields and options:

- GenBank IDs (one per line or comma separated) [optional]:** An empty text input field.
- Gene(s) [ex: 16S, COI] [optional]:** A text input field containing "COI" and a "Select genes" button to its right.
- Organism(s) [optional]:** A text input field containing "Anura".
- Publication term (title, authors, year) [optional]:** An empty text input field.
- Any additional query [ex: NOT sp.] [optional]:** An empty text input field.
- Refine your search terms to [optional]:** Four unchecked checkboxes: "Mitochondrial gene", "Mitogenome", "Chloroplast", and "Delete unannotated".
- At the bottom, there are two buttons: "◀ Previous" and "Next ▶".

212

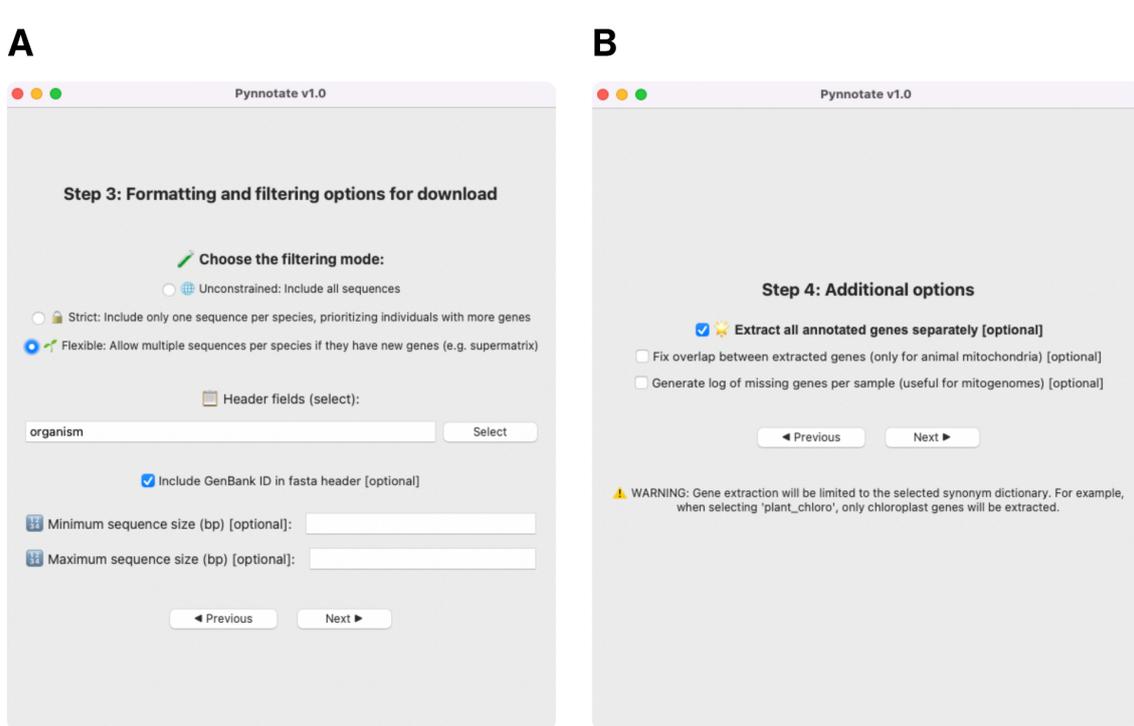
213 **Figure 2.** Pynnotate graphic interface. Step 2: Build a custom search query.

214

215 *c. Filtering and other options*

216 Besides sequence search, pynnotate offers many filtering options that can be applied to  
217 searches before results are returned to the user (Figure 3A). One of the most important features  
218 is the different filtering modes pynnotate allows. First, 'unconstrained' includes all available  
219 sequences regardless of redundancy, useful when users want to manually explore or curate all  
220 records. Second, 'strict' selects a single record per species by maximising the number of  
221 requested genes present in the annotation table. When multiple records contain the same  
222 number of target genes, pynnotate prioritises records with longer cumulative gene length and  
223 complete feature annotations. Ties are resolved deterministically to ensure reproducibility.  
224 Finally, 'flexible' selects one representative sequence per species for each gene independently.  
225 As a consequence, different genes within a multi-gene dataset may derive from distinct  
226 individuals of the same species. This approach maximizes gene occupancy across loci (e.g., in  
227 supermatrix analyses), but should be used with awareness of its implications for downstream  
228 phylogenetic concatenation. Both strict and flexible modes rely on the synonym dictionary to  
229 extract genes from each species. Furthermore, in Step 3, users can specify the header format  
230 for each sequence (e.g., organism name and accession number) and set minimum and  
231 maximum sequence lengths in base pairs.

232



233

234 **Figure 3.** Pynnotate graphic interface. **A.** Step 3: Formatting and filtering options for

235 download. **B.** Step 4: Additional options.

236

237 The additional options available in Step 4 control how results are presented (Figure

238 3B). The first particularly useful option is the extraction of all annotated genes separately,

239 grouping different organisms into respective files for each gene. Importantly, gene extraction

240 will again be limited by the selected synonym dictionary. The next two options allow

241 adjustment of the overlap between extracted mitochondrial genes and the generation of a log

242 documenting missing organisms per sample.

243

#### 244 *d. Download and output*

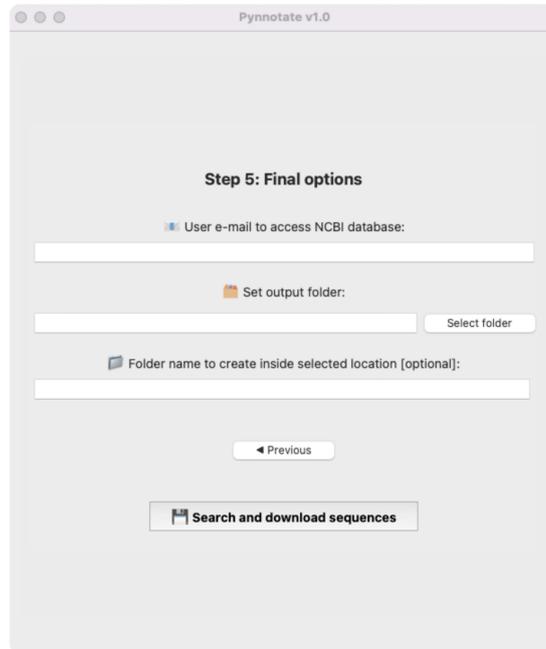
245 Finally, Step 5 includes saving and configuration options (Figure 4). First, users must

246 provide their email address to access the NCBI database. Then, users can select the output

247 folder and specify the folder name where pynnotate will save the results. Upon clicking on

248 ‘Search and download sequences’, several windows are displayed to provide informative

249 messages to users. When the download begins, a progress bar appears, indicating the download  
250 status. Finally, when the results are ready, an information window notifies the user.  
251



252  
253 **Figure 4.** Pynnotate graphic interface. Step 5: Final options.

254

#### 255 **4. Terminal workflow**

256 For users who prefer command-line execution or need to integrate pynnotate into  
257 automated pipelines, the tool also provides a terminal-based interface. The terminal version  
258 offers the same functionality as the graphical interface but uses a YAML configuration file  
259 (example available at  
260 <https://github.com/fernandacaron/pynnotate/blob/main/pynnotate/examples/config.yaml>) to  
261 define all parameters, making it suitable for batch processing and reproducible analyses.

262

##### 263 *a. Configuration file setup*

264 All execution parameters are specified in a YAML configuration file, which centralises  
265 the settings and eliminates the need for multiple command-line arguments. The configuration

266 file includes the same parameters available in the GUI: genome type, search terms or accession  
267 numbers, filtering options, output preferences, and file paths. An example configuration file is  
268 provided in the repository's examples folder.

269

#### 270 *b. Execution and output*

271 Once the configuration file is prepared, pynnotate can be executed with a single  
272 command:

273

```
bash pynnotate --config config.yaml
```

274

275 The terminal version generates the same output structure as the graphic interface.  
276 Progress information and any errors are displayed directly in the terminal, providing immediate  
277 feedback during execution. This approach is particularly valuable for high-throughput  
278 analyses, integration with other bioinformatics tools, and ensuring reproducibility through  
279 version-controlled configuration files.

280

## 281 **5. Discussion**

282 Pynnotate fills a gap in the phylogenetic annotation workflow by providing a simple  
283 yet flexible tool for extracting and organising gene sequences across multiple taxa. Pynnotate  
284 offers both an intuitive graphical interface and a fully scriptable command-line mode, making  
285 it accessible to users with or without programming experience. Because it is implemented as a  
286 modular Python package, it can be easily incorporated into larger bioinformatic pipelines or  
287 expanded by advanced users.

288 The main strengths of Pynnotate lie in its automated handling of repetitive annotation  
289 tasks and its ability to standardise output across heterogeneous GenBank records. Its focus on

290 batch annotation and feature retrieval makes it highly efficient for common preparatory steps  
291 in phylogenetic and comparative genomics studies. Future improvements may include broader  
292 support for additional feature types, identifying unannotated GenBank sequences through  
293 BLAST, and integration with alignment and phylogenetic inference tools. By combining a  
294 graphical interface for teaching environments with a command-line interface for reproducible  
295 research workflows, pynnotate bridges the gap between educational use and research-scale  
296 dataset assembly. This integration lowers the barrier to dataset preparation and has the potential  
297 to streamline early stages of evolutionary and molecular analyses.

298

## 299 6. References

- 300 Baykal, P. I., Łabaj, P. P., Markowetz, F., Schriml, L. M., Stekhoven, D. J., Mangul, S., &  
301 Beerenwinkel, N. (2024). Genomic reproducibility in the bioinformatics era. *Genome*  
302 *Biology*, 25(1), 213. <https://doi.org/10.1186/s13059-024-03343-2>
- 303 Borstein, S. R., & O'Meara, B. C. (2018). AnnotationBustR: an R package to extract  
304 subsequences from GenBank annotations. *PeerJ*, 6, e5179.  
305 <https://doi.org/10.7717/peerj.5179>
- 306 Cokelaer, T., Cohen-Boulakia, S., & Lemoine, F. (2023). Reprohackathons: Promoting  
307 reproducibility in bioinformatics through training. *Bioinformatics*, 39(Supplement\_1),  
308 i11–i20. <https://doi.org/10.1093/bioinformatics/btad227>
- 309 Emery, L. R., & Morgan, S. L. (2017). The application of project-based learning in  
310 bioinformatics training. *PLOS Computational Biology*, 13(8), e1005620.  
311 <https://doi.org/10.1371/journal.pcbi.1005620>
- 312 Hagen, J. B. (2011). The Origin and Early Reception of Sequence Databases. In M. Hamacher,  
313 M. Eisenacher, & C. Stephan (Eds.), *Data Mining in Proteomics* (Vol. 696, pp. 61–77).  
314 Humana Press. [https://doi.org/10.1007/978-1-60761-987-1\\_4](https://doi.org/10.1007/978-1-60761-987-1_4)

315 Harris, B. N., McCarthy, P. C., Wright, A. M., Schutz, H., Boersma, K. S., Shepherd, S. L.,  
316 Manning, L. A., Malisch, J. L., & Ellington, R. M. (2020). From panic to pedagogy:  
317 Using online active learning to promote inclusive instruction in ecology and  
318 evolutionary biology courses and beyond. *Ecology and Evolution*, *10*(22), 12581–  
319 12612. <https://doi.org/10.1002/ece3.6915>

320 Leray, M., Knowlton, N., Ho, S.-L., Nguyen, B. N., & Machida, R. J. (2019). GenBank is a  
321 reliable resource for 21st century biodiversity research. *Proceedings of the National*  
322 *Academy of Sciences*, *116*(45), 22651–22656.  
323 <https://doi.org/10.1073/pnas.1911714116>

324 Toelch, U., & Ostwald, D. (2018). Digital open science—Teaching digital tools for  
325 reproducible and transparent research. *PLOS Biology*, *16*(7), e2006022.  
326 <https://doi.org/10.1371/journal.pbio.2006022>

327