

Barcode Validator: A Python toolkit for structural and taxonomic validation of DNA barcode sequences

D. S. J. Groenengen ¹, Daniel A. J. Parsons   ², Ben Price   ², and Rutger A. Vos  

¹ Naturalis Biodiversity Center, Leiden, The Netherlands ² Natural History Museum, London, United Kingdom ¶ Corresponding author

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Open Journals](#) ↗

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

DNA barcoding has become a cornerstone technique in molecular biodiversity research, enabling rapid species identification and discovery through standardized genetic markers. The Barcode Validator is a Python toolkit designed to ensure the quality and accuracy of DNA barcode sequences before submission to public databases such as the Barcode of Life Data System (BOLD) and institutional repositories. The software performs both structural validation (assessing sequence quality, length, ambiguous bases, and marker-specific features like stop codons) and taxonomic validation (verifying specimen identifications through reverse taxonomy using BLAST-based identification services). Developed to support large-scale biodiversity genomics initiatives, particularly the Biodiversity Genomics Europe (BGE) and ARISE projects, the toolkit provides automated workflows for processing thousands of sequences with flexible configuration options and comprehensive reporting. The source code is available on GitHub at https://github.com/naturalis/barcode_validator under the Apache License 2.0, with distribution via PyPI and Bioconda, and a Galaxy tool wrapper for web-based access.

Statement of need

DNA barcoding projects generate large volumes of sequence data that must meet stringent quality standards before deposition in public databases. Manual validation of sequences is time-consuming, error-prone, and impractical for projects processing hundreds or thousands of specimens. Furthermore, the increasing adoption of genome skimming and high-throughput sequencing technologies produces multiple assembly attempts per specimen, requiring intelligent selection of the best valid sequence among alternatives.

The Barcode Validator addresses these needs by providing:

1. **Integrated validation:** Combined structural and taxonomic validation in a single workflow, with support for marker-specific requirements (e.g., stop codon detection for protein-coding genes via translation, GC content assessment for non-coding markers).
2. **Assembly triage:** Automatic selection of the best valid sequence when multiple assembly attempts exist per specimen, using configurable criteria including validation results and optional assembly quality metrics.
3. **Flexible taxonomic validation:** Support for multiple identification backends (BOLD API, local BLAST, Galaxy web services) and taxonomic backbones (BOLD, NCBI, Netherlands Species Register), enabling validation against expected specimen identifications at configurable taxonomic ranks.

39 4. **Batch processing:** Efficient handling of large datasets through batched API calls and
40 parallel processing where appropriate.

41 42. **Workflow integration:** Command-line interface suitable for automated pipelines, with
Galaxy tool integration for web-based access.

43 The software has been deployed in production workflows at Naturalis Biodiversity Center
44 (the Netherlands) and the Natural History Museum (United Kingdom) for the BGE project,
45 processing thousands of arthropod COI sequences from genome skimming experiments, and
46 the ARISE project for the validation of thousands of freshly sequenced vertebrate and marine
47 and terrestrial invertebrate specimens. Its design supports the quality assurance requirements
48 of modern DNA barcoding initiatives while remaining flexible enough to accommodate diverse
49 project-specific workflows.

50 State of the field

51 Existing tools for DNA barcode quality assurance are typically limited in scope. Quality control
52 tools like FastQC ([Andrews, 2010](#)) assess raw read quality rather than assembled barcode
53 sequences. Taxonomic identification tools like BOLD's identification engine ([Ratnasingham
& Hebert, 2007](#)) or standalone BLAST ([Altschul et al., 1990](#)) provide species identification
54 but lack integration with structural quality metrics. Biopython ([Cock et al., 2009](#)) provides
55 sequence manipulation and translation capabilities but not marker-specific HMM alignment for
56 reading frame detection. Profile HMM tools such as HMMER ([Eddy, 2011](#)) enable sequence
57 alignment but do not incorporate downstream validation logic.

58 No comprehensive solution exists that integrates structural validation, taxonomic verification,
59 and assembly triage for the specific workflow requirements of modern barcoding projects. The
60 Barcode Validator's contribution lies in this integration: combining HMM-based codon phase
61 detection, taxon-aware translation table selection, multi-backend taxonomic validation, and
62 assembly triage into a single configurable pipeline. No existing package provided extension
63 points suitable for adding this combined functionality; the unique combination of requirements
64 necessitated new software.

66 Software design

67 Barcode Validator is implemented in Python (3.9+) with a modular, extensible architecture
68 built around several key design patterns:

- 69 ▪ **Strategy pattern** for validators: An abstract Validator base class defines the validation
70 interface, with concrete implementations for structural validation (StructuralValidator
71 with subclasses ProteinCodingValidator and NonCodingValidator) and taxonomic
72 validation (TaxonomicValidator).
- 73 ▪ **Factory pattern** for services: Pluggable identification services (IDService hierarchy
74 supporting BOLD, BLAST, and Galaxy backends) and taxonomic resolvers
75 (TaxonResolver supporting BOLD, NCBI, and NSR taxonomies) enable flexible
76 backend selection.
- 77 ▪ **Orchestration pattern:** A ValidationOrchestrator coordinates the validation pipeline,
78 managing validator initialization, batch processing, result aggregation, and output
79 generation.

80 The central design decision was to separate validation logic (what measurements to collect)
81 from validity adjudication (what thresholds constitute pass/fail). This separation enables
82 the same codebase to serve diverse projects with different quality requirements—genome
83 skimming workflows demanding zero ambiguous bases versus Sanger sequencing tolerating
84 several—without code modifications. The Strategy pattern for validators and Factory pattern

85 for services enable runtime selection of validation approaches and identification backends. While
86 this introduces abstraction overhead, it proved essential for accommodating the consortium's
87 heterogeneous infrastructure: some partners operate local BLAST databases, others rely on
88 Galaxy web services, and still others use BOLD's identification API directly.

89 The software integrates with established bioinformatics tools including BLAST+ ([Camacho et al., 2009](#)) for sequence similarity searches, HMMER ([Eddy, 2011](#)) for profile Hidden Markov
90 Model-based alignment and codon phase detection, and Biopython ([Cock et al., 2009](#)) for
91 sequence manipulation and translation. External validation is performed through REST API
92 calls to BOLD ([Ratnasingham & Hebert, 2007](#)) and Galaxy ([The Galaxy Community, 2022](#))
93 identification services.

94 Input data can be provided as FASTA files with optional CSV metadata and BOLD Excel
95 spreadsheets containing specimen and taxonomic information. Validation results are output in
96 both human-readable TSV format (with detailed pass/fail status for each validation criterion)
97 and filtered FASTA format (containing only sequences meeting all validation requirements).

99 Research impact statement

100 The Barcode Validator has demonstrated substantial realized impact through its deployment in
101 the Biodiversity Genomics Europe (BGE) project. As documented in BGE Deliverable D8.4, the
102 toolkit processed sequences from over 18,500 specimens across 68 taxonomic orders, enabling
103 the submission of more than 47,000 validated DNA barcode sequences to BOLD and the
104 European Nucleotide Archive by October 2025. The validation framework identified systematic
105 issues including plate-swap errors that would otherwise have corrupted database submissions,
106 and revealed taxonomic patterns in validation success rates ranging from 0% to 100% across
107 orders—insights that directly informed protocol optimizations.

108 The software is deployed in production workflows at Naturalis Biodiversity Center (Netherlands)
109 and the Natural History Museum (United Kingdom), with the ARISE project using it for
110 validation of freshly sequenced vertebrate and invertebrate specimens. Community readiness is
111 evidenced by: distribution through PyPI and Bioconda channels; availability as a Galaxy tool
112 wrapper enabling web-based access for non-technical users; comprehensive documentation
113 including architecture diagrams and use-case examples; an Apache 2.0 license; and a public
114 GitHub repository with contribution guidelines. The toolkit's analytical outputs informed
115 the BGE consortium's understanding of genome skimming assembly parameter optimization,
116 demonstrating that the combination of specific preprocessing steps (`fcleaner=True`,
117 `merge=False`) with alignment thresholds ($r=1.0$, $s=50$) maximizes barcode recovery while
118 maintaining stringent quality standards.

119 AI usage disclosure

120 The overall software architecture—including the Strategy pattern for validators, Factory pattern
121 for services, and the separation of validation logic from criteria-based adjudication—was
122 conceived by the author prior to the widespread availability of usable large language models,
123 drawing on established object-oriented design principles. The parameterization of validation
124 logic, including marker-specific thresholds, taxonomic validation levels, and quality criteria, was
125 determined through iterative discussions among consortium users based on empirical analysis
126 of validation outcomes.

127 However, portions of the implementation benefited from generative AI assistance. Specifically,
128 Claude (Anthropic) and ChatGPT (OpenAI) were used to accelerate code syntax generation
129 for routine operations, produce initial drafts of docstrings and inline documentation, and refine
130 error handling patterns. The author reviewed, tested, and modified all AI-generated content

131 before incorporation. This manuscript was drafted by the author with AI assistance for prose
132 refinement and structural suggestions.

133 Acknowledgements

134 This work was supported by the Biodiversity Genomics Europe (BGE) project, which has received
135 funding from the European Union's Horizon Europe Research and Innovation Programme
136 under grant agreement No. 101059492, and the ARISE project. The author acknowledges the
137 Naturalis Biodiversity Center for institutional support and infrastructure, and Dick Groenenberg,
138 Dan Parsons and Ben Price for extensive testing, feedback, and improvement suggestions over
139 the course of this project.

140 References

- 141 Altschul, S. F., Gish, W., Miller, W., Myers, E. W., & Lipman, D. J. (1990). Basic local
142 alignment search tool. *Journal of Molecular Biology*, 215(3), 403–410. [https://doi.org/10.1016/S0022-2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2)
- 144 Andrews, S. (2010). *FastQC: A quality control tool for high throughput sequence data*.
145 <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
- 146 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., & Madden,
147 T. L. (2009). BLAST+: Architecture and applications. *BMC Bioinformatics*, 10(1), 421.
148 <https://doi.org/10.1186/1471-2105-10-421>
- 149 Cock, P. J., Antao, T., Chang, J. T., Chapman, B. A., Cox, C. J., Dalke, A., Friedberg,
150 I., Hamelryck, T., Kauff, F., Wilczynski, B., & Hoon, M. J. de. (2009). Biopython:
151 Freely available Python tools for computational molecular biology and bioinformatics.
152 *Bioinformatics*, 25(11), 1422–1423. <https://doi.org/10.1093/bioinformatics/btp163>
- 153 Eddy, S. R. (2011). Accelerated profile HMM searches. *PLoS Computational Biology*, 7(10),
154 e1002195. <https://doi.org/10.1371/journal.pcbi.1002195>
- 155 Ratnasingham, S., & Hebert, P. D. (2007). BOLD: The barcode of life data system. *Molecular
Ecology Notes*, 7(3), 355–364. <https://doi.org/10.1111/j.1471-8286.2007.01678.x>
- 157 The Galaxy Community. (2022). The Galaxy platform for accessible, reproducible and
158 collaborative biomedical analyses: 2022 update. *Nucleic Acids Research*, 50(W1),
159 W345–W351. <https://doi.org/10.1093/nar/gkac247>