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

**Rutger A. Vos** [1]¶, **D. S. J. Groenenberg** [1], **Daniel A. J. Parsons** [2], and **Ben Price** [2]

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

## 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.

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

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

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

# State of the field

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

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

# Software design

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

- **Strategy pattern** for validators: An abstract `Validator` base class defines the validation interface, with concrete implementations for structural validation (`StructuralValidator` with subclasses `ProteinCodingValidator` and `NonCodingValidator`) and taxonomic validation (`TaxonomicValidator`).

- **Factory pattern** for services: Pluggable identification services (`IDService` hierarchy supporting BOLD, BLAST, and Galaxy backends) and taxonomic resolvers (`TaxonResolver` supporting BOLD, NCBI, and NSR taxonomies) enable flexible backend selection.

- **Orchestration pattern**: A `ValidationOrchestrator` coordinates the validation pipeline, managing validator initialization, batch processing, result aggregation, and output generation.

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

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

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

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

## Research impact statement

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

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

## AI usage disclosure

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

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

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

## Acknowledgements

## References

Altschul, S. F., Gish, W., Miller, W., Myers, E. W., & Lipman, D. J. (1990). Basic local alignment search tool. *Journal of Molecular Biology*, *215*(3), 403–410. https://doi.org/10.1016/S0022-2836(05)80360-2

Andrews, S. (2010). *FastQC: A quality control tool for high throughput sequence data*. https://www.bioinformatics.babraham.ac.uk/projects/fastqc/

Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., & Madden, T. L. (2009). BLAST+: Architecture and applications. *BMC Bioinformatics*, *10*(1), 421. https://doi.org/10.1186/1471-2105-10-421

Cock, P. J., Antao, T., Chang, J. T., Chapman, B. A., Cox, C. J., Dalke, A., Friedberg, I., Hamelryck, T., Kauff, F., Wilczynski, B., & Hoon, M. J. de. (2009). Biopython: Freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics*, *25*(11), 1422–1423. https://doi.org/10.1093/bioinformatics/btp163

Eddy, S. R. (2011). Accelerated profile HMM searches. *PLoS Computational Biology*, *7*(10), e1002195. https://doi.org/10.1371/journal.pcbi.1002195

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

The Galaxy Community. (2022). The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2022 update. *Nucleic Acids Research*, *50*(W1), W345–W351. https://doi.org/10.1093/nar/gkac247