

# <sup>1</sup> BactScout: A Python pipeline for quality assessment and taxonomic profiling of bacterial sequencing data

<sup>3</sup> **Nabil-Fareed Alikhan**  <sup>1,3</sup>, **Varun Shamanna**<sup>2</sup>, and **GHRU Project Contributors**<sup>3</sup>

<sup>5</sup> 1 Centre for Genomic Pathogen Surveillance, University of Oxford, United Kingdom  
<sup>6</sup> 2 KIMS Global Healthcare and Research Unit

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Published: unpublished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](#))

## <sup>7</sup> Summary

BactScout is a Python-based pipeline for rapid, standardized quality assessment and taxonomic profiling of sequencing data from cultured bacterial isolates. It integrates tools like Fastp for read quality control, Sylph for species-level taxonomic profiling, and StringMLST for multi-locus sequence typing into a single, reproducible workflow. BactScout evaluates sequencing data across multiple quality dimensions—read quality, coverage depth, species purity, GC content, and strain typing—producing clear, interpretable quality metrics for downstream applications such as genome assembly, antimicrobial resistance prediction, genotyping, and phylogenetic inference.

The pipeline features a modular and extensible architecture with configurable quality thresholds, parallel sample processing, and detailed per-sample and batch-level reporting. A single command can process hundreds of samples, automatically generating summaries and visual outputs suitable for laboratory or high-performance computing environments.

BactScout emphasizes reproducibility and ease of use through deterministic, containerized environments managed by Pixi, ensuring consistent results across platforms. By combining quality control, taxonomic profiling, and strain typing in a unified, automated framework, BactScout reduces manual effort and improves standardization in bacterial genomics workflows.

## <sup>24</sup> Statement of need

Quality assessment of bacterial sequencing data is a critical and often under-standardized step in genomic analysis pipelines, particularly for applications requiring high-confidence genome assemblies. Common challenges include contamination, low sequencing yield, poor read quality, and variable fragment lengths. While existing tools can report these metrics, interpreting their biological relevance typically requires manual assessment that depends on the species and sequencing context. This leads to inconsistent quality decisions and reduced reproducibility across projects and laboratories.

BactScout addresses this gap by providing an automated, standardized workflow for assessing sequencing quality of bacterial isolates. It integrates established, fast-performing tools with relatively low memory and CPU requirements to evaluate read quality, taxonomic purity, and strain identity, applying clear pass/fail criteria based on configurable thresholds and species-specific quality benchmarks (as defined in [QualiBact](#)).

The pipeline is designed for typical isolate sequencing tasks encountered in public health surveillance and clinical microbiology, where rapid and reproducible decisions on sample quality are essential before downstream analyses such as genome assembly or resistance prediction. By formalizing interpretation and integrating species-aware thresholds, BactScout reduces

<sup>41</sup> subjective decision-making and improves consistency in bacterial genomics quality control. We  
<sup>42</sup> envision BactScout as an initial rapid screening step to identify high-quality samples suitable  
<sup>43</sup> for further analysis, working along side more comprehensive pipelines, which would explore the  
<sup>44</sup> genome assembly quality in greater depth.

## <sup>45</sup> BactScout Development

<sup>46</sup> BactScout is implemented in **Python ( 3.11)** and designed for ease of deployment and extensi-  
<sup>47</sup> bility.

<sup>48</sup> The pipeline provides a **command-line interface (CLI)** built with *Typer*, supporting both individ-  
<sup>49</sup> ual and batch processing modes. Dependencies—including **Fastp**, **Sylph**, and **StringMLST**—are  
<sup>50</sup> fully containerized and managed through **Pixi**, ensuring deterministic environments across  
<sup>51</sup> Linux and macOS systems.

<sup>52</sup> System-level tests, mock data, and example configurations are included to validate installations  
<sup>53</sup> and future development. Continuous integration is performed via GitHub Actions on multiple  
<sup>54</sup> platforms (Ubuntu 22.04 and macOS latest).

<sup>55</sup> Parallelization is implemented through Python's thread pool executor, enabling efficient  
<sup>56</sup> processing of large sample sets on high-performance or cloud computing environments.

## <sup>57</sup> Tools Utilized in BactScout

<sup>58</sup> BactScout orchestrates three primary external tools to evaluate sequencing quality and taxo-  
<sup>59</sup> nomic composition (Table 1).

Tool	Function	Quality Dimension
<b>Fastp</b>	Read-level quality control and adapter trimming	Read quality
<b>Sylph</b>	Taxonomic profiling and species purity estimation	Species identification
<b>StringMLST</b>	Multi-locus sequence typing (MLST) assignment	Strain typing

<sup>60</sup> Each module outputs standardized JSON or tabular results that are parsed and evaluated  
<sup>61</sup> against BactScout's threshold schema. Quality decisions (PASS/WARNING/FAIL) are derived  
<sup>62</sup> from metrics such as: - Mean Q30 score and read length (**Fastp**) - Genome coverage and  
<sup>63</sup> species composition (**Sylph**) - MLST type validity and completeness (**StringMLST**)

<sup>64</sup> Default thresholds are defined in YAML configuration files and can be customized to project  
<sup>65</sup> or organism-specific standards.

## <sup>66</sup> Quality Assessment and Reporting

<sup>67</sup> BactScout performs quality evaluation in four primary domains:

- <sup>68</sup> **1. Read Quality:** Calculates mean read length and percentage of bases ( Q30) from **Fastp**  
<sup>69</sup> outputs.
- <sup>71</sup> **2. Coverage Depth:** Estimates genome coverage both from read counts and **Sylph**-derived  
<sup>72</sup> genome size.

74        3. **Species Purity:** Quantifies dominant species proportion and flags cross-species  
75            contamination.

76  
77        4. **Strain Typing:** Runs StringMLST when a single dominant species is detected to validate  
78            strain-level assignment.

79        Each sample is assigned an overall **status**—PASS, WARNING, or FAIL—with explanatory notes  
80            for each metric.

81        Reports include:  
82            - **Per-sample CSV summaries** with metric breakdowns  
83            - **Batch-level summaries** aggregating performance across all samples  
83            - **Optional Fastp HTML reports** for visual inspection

84        Outputs are human-readable and machine-parseable, facilitating integration with LIMS systems  
85            or downstream pipelines such as **Nextflow**.

## 86        Applications

87        BactScout is applicable across multiple domains:

- 88            ▪ **Genome assembly projects** – Pre-assembly QC to identify high-quality inputs
- 89
- 90            ▪ **Epidemiological surveillance** – Rapid strain verification and contamination detection
- 91
- 92            ▪ **Sequencing QA/QC** – Standardized acceptance criteria for clinical or public health  
93            laboratories
- 94
- 95            ▪ **Multi-center research cohorts** – Harmonized quality reporting across institutions

96        By producing interpretable, standardized outputs, BactScout helps ensure that only high-quality,  
97            biologically relevant sequencing data progress to downstream analysis.

## 98        Source Code and Documentation

99        Source code for **BactScout** is available at <https://github.com/ghruproject/bactscout> under  
100          the **GPLv3 License**.

101        Comprehensive documentation—covering installation, usage, configuration, and troubleshooting—  
102          is hosted on GitHub Pages.

103        It includes API references, example datasets, and developer contribution guidelines.

## 104        Acknowledgements

105        BactScout builds upon three outstanding open-source tools: **Fastp** (Chen et al., 2018), **Sylph**  
106          (Unckless et al., 2023), and **StringMLST** (Datta et al., 2016).

107        We thank contributors from the **Global Health Research Unit (GHRU)** for feedback during  
108          design and testing, and the open-source community for providing the foundational libraries  
109          and infrastructure enabling this work.

## 110        References