

¹ Tataki: Enhancing the robustness of bioinformatics workflows with simple, tolerant file format detection

³ **Masaki Fukui**  ¹, **Hirotaka Suetake**  ¹, and **Tazro Ohta**  ^{2,3}

⁴ 1 Sator, Inc. 2 Institute for Advanced Academic Research, Chiba University  3 Department of Artificial Intelligence Medicine, Graduate School of Medicine, Chiba University

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

Software

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

Editor: 

Submitted: 25 December 2025

Published: unpublished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](https://creativecommons.org/licenses/by/4.0/)).

⁶ Summary

⁷ Tataki is a lightweight command-line tool for workflow-oriented file format detection in
⁸ bioinformatics. Automated genomics pipelines frequently exchange large intermediate files (e.g.,
⁹ SAM/BAM/CRAM, VCF, BED), which can sometimes be empty, truncated, or structurally
¹⁰ inconsistent. Because many downstream tools do not reliably detect such conditions, these
¹¹ malformed files can lead to silent failures that are hard to diagnose. Tataki inspects actual file
¹² contents using strict, domain-aware parsers to identify malformed or mixed-format files that
¹³ general-purpose tools often misclassify, and includes an extension mode based on the Common
¹⁴ Workflow Language (CWL) to support custom or emerging formats. It emphasizes verifying
¹⁵ that each file matches the expected format at its workflow stage, aiming to improve workflow
¹⁶ robustness and provenance.

¹⁷ Statement of need

¹⁸ Modern bioinformatics workflows integrate many specialized analytical tools to process large
¹⁹ scale sequencing data (Perkel, 2019). These workflows are designed to reduce manual
²⁰ intervention and to ensure reproducibility by automating complex multi-step analyses. However,
²¹ despite their widespread adoption, workflow executions often fail for surprisingly simple reasons:
²² intermediate files may be empty, truncated, or internally inconsistent, and many bioinformatics
²³ tools do not reliably signal such failures through exit codes (Niu et al., 2022). As a consequence,
²⁴ downstream tools may misinterpret file formats, propagate incorrect assumptions, or terminate
²⁵ unexpectedly. These silent errors reduce the overall robustness and fault tolerance of automated
²⁶ workflows, which becomes increasingly problematic as datasets grow and manual oversight
²⁷ becomes infeasible.

²⁸ This fragility is rooted in longstanding characteristics of the bioinformatics ecosystem. The
²⁹ field contains a large and heterogeneous collection of file formats, many of which lack formal
³⁰ specifications or have multiple variants used by different tools. Beyond a handful of well-
³¹ standardized formats such as SAM/BAM (Li et al., 2009)/CRAM (Cochrane et al., 2012),
³² VCF (Danecek et al., 2011) or BED (Niu et al., 2022), many commonly used formats have
³³ ambiguous boundaries or are interpreted differently across implementations (Rehm et al.,
³⁴ 2021). At the same time, bioinformatics software is often developed by individual research
³⁵ groups for specific tasks, typically originating as stand-alone research code created under tight
³⁶ time and resource constraints, leading to inconsistent behaviors, uneven error handling, and
³⁷ non-standard assumptions about input data (Brack et al., 2022). These factors combine to
³⁸ create a landscape in which workflows must routinely exchange files whose correctness cannot
³⁹ be assumed.

⁴⁰ To address these challenges, we developed Tataki, a lightweight and workflow-friendly file
⁴¹ format detection tool tailored to bioinformatics. Tataki is designed to combine strict parsers

42 with an extensive mode for structurally variable formats, aiming to support both standarized and
43 real-world data irregularities. By incorporating Tataki between workflow steps, researchers and
44 developers can rapidly identify format-related anomalies before they propagate, and improve
45 debugging of multi-step pipelines. Ultimately, this leads to more robust and reproducible
46 research, and enables more accurate provenance recording for workflow outputs.

47 Overview of Tataki

48 Tataki is written in Rust and designed as a workflow component rather than an interactive
49 validator. It examines file content directly without relying on filename extensions and handles
50 compressed inputs in gzip and bzip2 formats by decompressing them before analysis. To
51 detect anomalous or truncated files that superficially appear valid, Tataki scans entire files.
52 However, given that genomics intermediate files can reach hundreds of gigabytes, it also
53 provides bounded inspection (e.g., checking only the first N records or lines) to reduce latency.
54 Results are presented in a structured manner using terms from the EDAM ontology ([Black et al., 2022](#)).

55 Tataki supports two complementary detection modes. Its native mode uses strict, domain-
56 aware parsers for major genomics formats. For greater flexibility, the External Extension Mode
57 allows users to define format identification logic via CWL ([Crusoe et al., 2022](#)).

59 External Extension Mode

60 The External Extension Mode enables users to define custom file format identification logic
61 through the CWL, which adds support for emerging or project-specific formats without
62 modifying the core software. In this mode, format recognition is delegated to a user-supplied
63 CWL document, which specifies how a file should be processed and which EDAM format
64 identifier should be assigned upon successful validation.

65 External Extension Mode allows Tataki to remain lightweight while accommodating the diverse
66 and evolving landscape of file formats used in bioinformatics research. This approach provides a
67 flexible pathway for integrating domain-specific validators and facilitates more reliable workflow
68 execution in specialized or rapidly developing research areas.

69 State of the field

70 Existing file format detection tools only partially address the challenges of robust file validation in
71 genomics workflows. General-purpose utilities such as file ([Darwin, n.d.](#)), Magika ([Fratantonio et al., 2025](#)), Siegfried ([Lehane, n.d.](#)), and TrID ([Pontello, n.d.](#)) rely on static magic byte
72 rules, trained classifiers, heuristics, or other techniques, and have minimal or no support for
73 domain-specific bioinformatics formats. Due to these algorithmic characteristics, they often
74 misclassify genomics files with mixed or mislabeled content, typically identifying only the first
75 recognizable format segment.

77 PipeVal ([Patel et al., 2024](#)) is a tool developed specifically for use in bioinformatics workflows. It
78 performs quick validations using format-specific modules selected based on file extensions, and
79 includes checksum-based comparisons to detect file corruption. However, its scope differs from
80 Tataki: PipeVal assumes files are correctly typed and verifies internal consistency, whereas
81 Tataki focuses on independently verifying that the output format truly matches what is
82 expected at that workflow stage, even in the presence of malformed or hybrid files.

83 Given this landscape, extending existing tools would not address the structural limitations in
84 format coverage, content verification, or pipeline integration.

85 Software Design

86 Tataki was designed with an emphasis on execution flexibility and extensibility in workflow
87 environments. Extensibility here refers both to the ability for users to introduce project-specific
88 format checks and to keeping the codebase approachable for community contributions.

89 To support flexibility and extensibility, format detectors are implemented as independent
90 modules, each responsible for a single file format. This structure allows detectors to be
91 composed and reordered without affecting the core logic. In principle, such a design could
92 enable dynamic loading of user-developed detectors, reducing the size of the core binary and
93 allowing third-party extensions to be deployed without recompilation.

94 However, dynamically loading in Rust requires Rust's unsafe features, and introduces additional
95 constraints on version compatibility, testing, and deployment. In heterogeneous workflow
96 environments, these factors undermine reproducibility and complicate continuous integration
97 and distribution.

98 Instead of relying on native dynamic loading, Tataki adopts External Extension Mode as its
99 extensibility mechanism. This design balances flexibility with robustness, aligning Tataki's
100 extensibility model with the reliability requirements of automated bioinformatics workflows.

101 Research Impact Statement

102 Tataki improves the robustness of bioinformatics workflows by detecting malformed, truncated,
103 or structurally inconsistent intermediate files before they propagate to downstream analysis
104 steps.

105 Tataki is under consideration for integration into an existing peer-reviewed workflow execution
106 service developed within our organization. In this setting, Tataki would act as a shared preflight
107 layer, enabling format checks and machine-readable provenance without requiring users to
108 modify their workflow definitions. This provides benefits such as earlier failure detection,
109 reduced waste of computational resources, and more reliable provenance records across diverse
110 workflow engines.

111 Tataki is designed to be community-ready. Format detectors are implemented as independent
112 modules, and the repository provides documentation, templates, and tests to facilitate the
113 addition of new detectors. The project is released under the Apache-2.0 license. For ease
114 of adoption, Tataki is distributed as a single static binary and as an OCI container image,
115 allowing straightforward integration into local environments and container-based workflow
116 platforms.

117 Limitations

118 Tataki focuses on identifying file formats and detecting structural anomalies; it does not
119 perform semantic validation of biological content. For example, it does not verify whether
120 sequence identifiers are consistent across related files or whether genomic coordinates fall
121 within valid ranges. Such checks remain the responsibility of downstream analysis tools or
122 dedicated validation software.

123 AI Usage Disclosure

124 Generative AI tools were used in a limited manner during the development of Tataki and the
125 preparation of this manuscript. During software development, generative AI was used to assess
126 the feasibility of design ideas, clarify specific implementation approaches, and generate code
127 snippets. No agentic or autonomous AI systems were used for coding. The codebase was

¹²⁸ primarily written, reviewed, and validated by MF, with minor contributions from other human
¹²⁹ collaborators.

¹³⁰ Project documentation was primarily written by MF, with help from generative AI for grammar
¹³¹ checks and clarification of wording and tone. For manuscript preparation, generative AI was
¹³² used to assist with drafting text. All AI-assisted content was subsequently reviewed, edited,
¹³³ and verified by the authors to ensure technical accuracy and consistency with the software
¹³⁴ implementation and cited literature.

¹³⁵ Acknowledgements

¹³⁶ This work was supported by the research grant from the SECOM Science and Technology
¹³⁷ Foundation (FY2023 grant program).

¹³⁸ References

- ¹³⁹ Black, M., Lamothe, L., Eldakroury, H., Kierkegaard, M., Priya, A., Machinda, A., Singh
¹⁴⁰ Khanduja, U., Patoliya, D., Rathi, R., Che Nico, T. P., Umutesi, G., Blankenburg, C.,
¹⁴¹ Op, A., Chieke, P., Babatunde, O., Laurie, S., Neumann, S., Schwämmele, V., Kuzmin,
¹⁴² I., ... Kalaš, M. (2022). *EDAM: The bioscientific data analysis ontology (update 2021)*.
¹⁴³ <https://doi.org/10.7490/f1000research.1118900.1>
- ¹⁴⁴ Brack, P., Crowther, P., Soiland-Reyes, S., Owen, S., Lowe, D., Williams, A. R., Groom, Q.,
¹⁴⁵ Dillen, M., Coppens, F., Grüning, B., Eguinoia, I., Ewels, P., & Goble, C. (2022). Ten
¹⁴⁶ simple rules for making a software tool workflow-ready. *PLOS Computational Biology*,
¹⁴⁷ 18(3), 1–11. <https://doi.org/10.1371/journal.pcbi.1009823>
- ¹⁴⁸ Cochrane, G., Alako, B., Amid, C., Bower, L., Cerdeño-Tárraga, A., Cleland, I., Gibson, R.,
¹⁴⁹ Goodgame, N., Jang, M., Kay, S., Leinonen, R., Lin, X., Lopez, R., McWilliam, H., Oisel,
¹⁵⁰ A., Pakseresht, N., Pallreddy, S., Park, Y., Plaister, S., ... Zalunin, V. (2012). Facing
¹⁵¹ growth in the European Nucleotide Archive. *Nucleic Acids Research*, 41(D1), D30–D35.
¹⁵² <https://doi.org/10.1093/nar/gks1175>
- ¹⁵³ Crusoe, M. R., Abeln, S., Iosup, A., Amstutz, P., Chilton, J., Tijanić, N., Ménager, H., Soiland-
¹⁵⁴ Reyes, S., Gavrilović, B., Goble, C., & Community, T. C. (2022). Methods included:
¹⁵⁵ Standardizing computational reuse and portability with the Common Workflow Language.
¹⁵⁶ *Communications of the ACM*, 65(6), 54–63. <https://doi.org/10.1145/3486897>
- ¹⁵⁷ Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., Handsaker,
¹⁵⁸ R. E., Lunter, G., Marth, G. T., Sherry, S. T., McVean, G., Durbin, R., & Group, 1000. G.
¹⁵⁹ P. A. (2011). The variant call format and VCFtools. *Bioinformatics*, 27(15), 2156–2158.
¹⁶⁰ <https://doi.org/10.1093/bioinformatics/btr330>
- ¹⁶¹ Darwin, I. (n.d.). *Fine Free File Command (and libmagic) — darwinsys.com*. Retrieved
¹⁶² February 6, 2026, from <https://www.darwinsys.com/file/>
- ¹⁶³ Fratantonio, Y., Invernizzi, L., Farah, L., Thomas, K., Zhang, M., Albertini, A., Galilee, F.,
¹⁶⁴ Metitieri, G., Cretin, J., Petit-Bianco, A., Tao, D., & Bursztein, E. (2025). MAGIKA:
¹⁶⁵ AI-Powered Content-Type Detection. *2025 IEEE/ACM 47th International Conference on*
¹⁶⁶ *Software Engineering (ICSE)*, 2638–2649. <https://doi.org/10.1109/ICSE55347.2025.00158>
- ¹⁶⁷ Lehane, R. (n.d.). *IT for archivists — itforarchivists.com*. Retrieved February 6, 2026, from
¹⁶⁸ <https://www.itforarchivists.com/siegfried>
- ¹⁶⁹ Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis,
¹⁷⁰ G., Durbin, R., & Subgroup, 1000. G. P. D. P. (2009). The Sequence Alignment/Map
¹⁷¹ format and SAMtools. *Bioinformatics*, 25(16), 2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>

- 173 Niu, Y. N., Roberts, E. G., Denisko, D., & Hoffman, M. M. (2022). Assessing and assuring
174 interoperability of a genomics file format. *Bioinformatics*, 38(13), 3327–3336. <https://doi.org/10.1093/bioinformatics/btac327>
- 175
- 176 Patel, Y., Beshlikyan, A., Jordan, M., Kim, G., Holmes, A., Yamaguchi, T. N., & Boutros, P.
177 C. (2024). PipeVal: Light-weight extensible tool for file validation. *Bioinformatics*, 40(2),
178 btae079. <https://doi.org/10.1093/bioinformatics/btae079>
- 179 Perkel, J. M. (2019). Workflow systems turn raw data into scientific knowledge. *Nature*,
180 573(7772), 149–150. <https://doi.org/10.1038/d41586-019-02619-z>
- 181 Pontello, M. (n.d.). *Marco Pontello's Home - Software - TrID — mark0.net*. Retrieved
182 February 6, 2026, from <https://mark0.net/soft-trid-e.html>
- 183 Rehm, H. L., Page, A. J. H., Smith, L., Adams, J. B., Alterovitz, G., Babb, L. J., Barkley,
184 M. P., Baudis, M., Beauvais, M. J. S., Beck, T., Beckmann, J. S., Beltran, S., Bernick,
185 D., Bernier, A., Bonfield, J. K., Boughtwood, T. F., Bourque, G., Bowers, S. R., Brookes,
186 A. J., ... Birney, E. (2021). GA4GH: International policies and standards for data sharing
187 across genomic research and healthcare. *Cell Genomics*, 1(2), 100029. <https://doi.org/10.1016/j.xgen.2021.100029>
- 188