


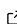


1 sourmash v4: A multitool to quickly search, compare,
2 and analyze genomic and metagenomic data sets

3 Luiz Irber ^{1*}, N. Tessa Pierce-Ward ^{1*}, Mohamed Abuelanin ¹, Harriet
4 Alexander ², Abhishek Anant ⁹, Keya Barve ¹, Colton Baumler ¹, Olga
5 Botvinnik ³, Phillip Brooks ¹, Daniel Dsouza ⁹, Laurent Gautier ⁹,
6 Mahmudur Rahman Hera ⁴, Hannah Eve Houts ¹, Lisa K. Johnson ¹,
7 Fabian Klötzl ⁵, David Koslicki ⁴, Marisa Lim ¹, Ricky Lim ⁹, Bradley
8 Nelson ⁹, Ivan Ogasawara ⁹, Taylor Reiter ¹, Camille Scott ¹, Andreas
9 Sjödin ⁶, Daniel Standage ⁷, S. Joshua Swamidass ⁸, Connor Tiffany ⁹,
10 Pranathi Vemuri ³, Erik Young ¹, and C. Titus Brown ^{1¶}

11 1 University of California, Davis 2 Woods Hole Oceanographic Institution 3 Chan-Zuckerberg Biohub 4
12 Pennsylvania State University 5 MPI for Evolutionary Biology 6 Swedish Defence Research Agency (FOI)
13 7 National Bioforensic Analysis Center 8 Washington University in St Louis 9 No affiliation ¶
14 Corresponding author * These authors contributed equally.

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Published: unpublished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](#))

15 Summary

16 sourmash is a command line tool and Python library for sketching collections of DNA, RNA,
17 and amino acid k-mers for biological sequence search, comparison, and analysis ([Pierce et al.,
18 2019](#)). sourmash's FracMinHash sketching supports fast and accurate sequence comparisons
19 between datasets of different sizes ([Irber, Brooks, et al., 2022](#)), including taxonomic profiling
20 ([Portik et al., 2022](#)), functional profiling ([Rahman Hera, Liu, et al., 2023](#)), and petabase-scale
21 sequence search ([Irber, Pierce-Ward, et al., 2022](#)). From release 4.x, sourmash is built on top
22 of Rust and provides an experimental Rust interface.

23 FracMinHash sketching is a lossy compression approach that represents data sets using a
24 "fractional" sketch containing $1/S$ of the original k-mers. Like other sequence sketching
25 techniques (e.g. MinHash, ([Ondov et al., 2015](#))), FracMinHash provides a lightweight way to
26 store representations of large DNA or RNA sequence collections for comparison and search.
27 Sketches can be used to identify samples, find similar samples, identify data sets with shared
28 sequences, and build phylogenetic trees. FracMinHash sketching supports estimation of overlap,
29 bidirectional containment, and Jaccard similarity between data sets and is accurate even for
30 data sets of very different sizes.

31 Since sourmash v1 was released in 2016 ([Brown & Irber, 2016](#)), sourmash has expanded to
32 support new database types and many more command line functions. In particular, sourmash
33 now has robust support for both Jaccard similarity and Containment calculations, which
34 enables analysis and comparison of data sets of different sizes, including large metagenomic
35 samples. As of v4.4, sourmash can convert these to estimated Average Nucleotide Identity
36 (ANI) values, which can provide improved biological context to sketch comparisons ([Rahman
37 Hera, Pierce-Ward, et al., 2023](#)).

38 Statement of Need

39 Large collections of genomes, transcriptomes, and raw sequencing data sets are readily
40 available in biology, and the field needs lightweight computational methods for searching and

41 summarizing the content of both public and private collections. sourmash provides a flexible set
42 of programmatic tools for this purpose, together with a robust and well-tested command-line
43 interface. It has been used in over 350 publications (based on citations of Brown & Irber
44 (2016) and Pierce et al. (2019)) and it continues to expand in functionality.

45 Acknowledgements

46 This work was funded in part by the Gordon and Betty Moore Foundation's Data-Driven
47 Discovery Initiative [GBMF4551 to CTB]. It is also funded in part by the National Science
48 Foundation [#2018522 to CTB] and PIG-PARADIGM (Preventing Infection in the Gut of
49 developing Piglets—and thus Antimicrobial Resistance – by disentangling the interface of Dlet,
50 the host and the Gastrointestinal Microbiome) from the Novo Nordisk Foundation to CTB.

51 Notice: This manuscript has been authored by BNBI under Contract No. HSHQDC-15-C-00064
52 with the DHS. The US Government retains and the publisher, by accepting the article for
53 publication, acknowledges that the USG retains a non-exclusive, paid-up, irrevocable, world-
54 wide license to publish or reproduce the published form of this manuscript, or allow others to
55 do so, for USG purposes. Views and conclusions contained herein are those of the authors and
56 should not be interpreted to represent policies, expressed or implied, of the DHS.

57 References

- 58 Brown, C. T., & Irber, L. (2016). Sourmash: A library for MinHash sketching of DNA. *Journal*
59 *of Open Source Software*, 1(5), 27. <https://doi.org/10.21105/joss.00027>
- 60 Irber, L. C., Brooks, P. T., Reiter, T. E., Pierce-Ward, N. T., Hera, M. R., Koslicki, D., & Brown,
61 C. T. (2022). Lightweight compositional analysis of metagenomes with FracMinHash and
62 minimum metagenome covers. *bioRxiv*. <https://doi.org/10.1101/2022.01.11.475838>
- 63 Irber, L. C., Pierce-Ward, N. T., & Brown, C. T. (2022). Sourmash branchwater enables
64 lightweight petabyte-scale sequence search. *bioRxiv*. <https://doi.org/10.1101/2022.11.02.514947>
- 65
- 66 Ondov, B. D., Treangen, T. J., Mallonee, A. B., Bergman, N. H., Koren, S., & Phillippy, A.
67 M. (2015). Fast genome and metagenome distance estimation using MinHash. *bioRxiv*,
68 029827. <https://doi.org/10.1101/029827>
- 69 Pierce, N. T., Irber, L., Reiter, T., Brooks, P., & Brown, C. T. (2019). Large-scale se-
70 quence comparisons with sourmash. *F1000Research*, 8, 1006. <https://doi.org/10.12688/f1000research.19675.1>
- 71
- 72 Portik, D. M., Brown, C. T., & Pierce-Ward, N. T. (2022). Evaluation of taxonomic
73 profiling methods for long-read shotgun metagenomic sequencing datasets. *Bioinformatics*.
74 <https://doi.org/10.1186/s12859-022-05103-0>
- 75 Rahman Hera, M., Liu, S., Wei, W., Rodriguez, J. S., Ma, C., & Koslicki, D. (2023). Fast,
76 lightweight, and accurate metagenomic functional profiling using FracMinHash sketches.
77 *bioRxiv*, 2023–2011.
- 78 Rahman Hera, M., Pierce-Ward, N. T., & Koslicki, D. (2023). Deriving confidence intervals for
79 mutation rates across a wide range of evolutionary distances using FracMinHash. *Genome*
80 *Research*, gr-277651. <https://doi.org/10.1101/gr.277651.123>