

# <sup>1</sup> KrakenParser: Efficient Post-processing of <sup>2</sup> Kraken2-like Taxonomic Reports

<sup>3</sup> Ilia V. Popov  

<sup>4</sup> 1 Faculty of Bioengineering and Veterinary Medicine, Don State Technical University, Russia 

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

## Software

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

Editor: Claudia Solis-Lemus  

## Reviewers:

- [@johanneswerner](#)

Submitted: 01 June 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Published: unpublished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](#)).

## <sup>5</sup> Summary

<sup>6</sup> KrakenParser is an open-source software tool (with a command-line interface and Python API) designed to streamline the post-analysis of metagenomic classification results produced by <sup>7</sup> Kraken2 ([Wood et al., 2019](#)) and similar taxonomic profilers such as Bracken ([Lu et al., 2017](#)) <sup>8</sup> and Metabuli ([J. Kim & Steinegger, 2024](#)). Kraken2 is a widely used taxonomic classifier <sup>9</sup> that assigns metagenomic reads to taxa using exact k-mer matches, achieving high speed and <sup>10</sup> accuracy. However, the raw output of Kraken2 ([Wood et al., 2019](#)) (and related tools) is <sup>11</sup> a text report that can be cumbersome to interpret and aggregate across multiple samples. <sup>12</sup> KrakenParser addresses this need by converting multiple Kraken-format reports into structured <sup>13</sup> tables (CSV files) at various taxonomic ranks (from phylum down to species), performing <sup>14</sup> filtering and normalization (including relative abundance calculations), and providing APIs to <sup>15</sup> produce publication-ready plots. The tool automates the multi-step process of combining and <sup>16</sup> cleaning Kraken results, allowing researchers to quickly obtain human-readable summaries of <sup>17</sup> community composition. KrakenParser's focus is on efficiency, ease-of-use, and integration: <sup>18</sup> it can run an entire conversion pipeline with a single command and also be imported as a <sup>19</sup> Python library for custom workflows. In summary, KrakenParser significantly reduces the <sup>20</sup> manual effort required to post-process metagenomic classification data, enabling scientists to <sup>21</sup> go from raw classifier output to analysis-ready tables and figures in one step.

## <sup>23</sup> Statement of need

<sup>24</sup> Analyzing the taxonomic profiles of metagenomic samples often involves running k-mer based <sup>25</sup> classifiers (like Kraken2) that generate detailed reports of read counts and abundances across <sup>26</sup> taxa. These reports, while information-rich, are not immediately convenient for comparative <sup>27</sup> analysis: they list each taxon in a hierarchical format for a single sample, and researchers must <sup>28</sup> manually parse and merge multiple files to compare communities across samples. Existing <sup>29</sup> scripts such as the KrakenTools suite ([Lu et al., 2022](#)) (developed alongside Kraken) provide <sup>30</sup> some post-processing functionality, but they require multiple steps and technical expertise to <sup>31</sup> use. Similarly, interactive tools like Pavian focus on visualization and exploration of Kraken <sup>32</sup> results rather than automated batch processing ([Breitwieser & Salzberg, 2020](#)). There is a <sup>33</sup> clear need for a streamlined solution to transform raw Kraken-family outputs into tidy data <sup>34</sup> matrices and summary statistics that can be readily used in downstream analysis or publication <sup>35</sup> figures. KrakenParser fulfills this need by offering an all-in-one pipeline that reads in multiple <sup>36</sup> Kraken2/Bracken/Metabuli reports and outputs clean CSV tables of taxonomic counts or <sup>37</sup> relative abundances, optionally filtering out low-abundance taxa or non-target taxa (e.g. human <sup>38</sup> reads) as specified by the user. This greatly simplifies metagenomic workflows, especially in <sup>39</sup> comparative studies or clinical settings where dozens of samples must be processed consistently. <sup>40</sup> By bridging the gap between raw classifier output and statistical analysis, KrakenParser <sup>41</sup> empowers researchers who may not be bioinformatics experts to leverage high-throughput <sup>42</sup> metagenomics with minimal data wrangling.

43 Metagenomic classification has seen rapid development, with numerous tools available for  
44 assigning sequencing reads to taxa. Kraken was introduced in 2014 as an ultrafast k-mer based  
45 classifier ([Wood & Salzberg, 2014](#)), and its successor Kraken2 ([Wood et al., 2019](#)) further  
46 reduced memory usage and improved speed . Other k-mer classifiers include Bracken ([Lu et](#)  
47 [al., 2017](#)), which refines Kraken's counts to improve abundance estimates, KrakenUniq which  
48 tracks unique k-mers per taxon to reduce false positives ([Breitwieser et al., 2018](#)), Centrifuge  
49 which uses an FM-index to allow classification with compressed databases ([D. Kim et al.,](#)  
50 [2016](#)), and CLARK which uses discriminative k-mers for fast classification ([Ounit et al., 2015](#)).  
51 More recently, tools like Kaiju perform classification in protein space for greater sensitivity  
52 (especially on viruses) ([Menzel et al., 2016](#)), and Metaboli combines DNA and translated  
53 amino acid matching to improve accuracy ([J. Kim & Steinegger, 2024](#)). Comprehensive  
54 evaluations have benchmarked these methods' accuracy and speed, and community challenges  
55 like CAMI have pushed development of improved classifiers ([Sczyrba et al., 2017](#)). Despite the  
56 variety of classifiers, a common challenge remains: the output format. Many tools output  
57 reports similar to Kraken's: tab-delimited text with hierarchical labels and counts. To interpret  
58 such outputs, researchers often rely on additional scripts or manual processing. KrakenTools  
59 ([Lu et al., 2022](#)) provides scripts to combine Kraken reports, convert to other formats (e.g.,  
60 Krona for visualization). Pavian and other interactive platforms allow users to visualize results  
61 with Sankey diagrams and heatmaps ([Breitwieser & Salzberg, 2020](#)), but require use of a  
62 web interface or R environment. There are also lightweight utilities (e.g., [spideog](#)) to convert  
63 Kraken reports to CSV or clean them, and researchers adept in programming sometimes write  
64 custom parsing scripts. In summary, prior to KrakenParser, users had to piece together multiple  
65 tools to achieve tasks like merging reports from multiple samples, summing reads at specific  
66 taxonomic ranks, and computing relative abundances. KrakenParser builds on this state of  
67 the field by consolidating the post-processing steps into one tool. It serves as an ideological  
68 successor to KrakenTools ([Lu et al., 2022](#)), using some of the same internal conversion  
69 steps (like KrakenTools' report-to-MPA conversion) but adding improvements in automation,  
70 filtering, and output formatting. By producing standardized CSV tables (with samples as rows  
71 and taxa as columns) and by computing percentages automatically, KrakenParser greatly  
72 accelerates the transition from raw classification data to biological insights. This is particularly  
73 valuable given the increasing scale of metagenomic studies (where dozens or hundreds of  
74 samples are profiled) and the need for reproducible, efficient analysis pipelines.

## 75 Implementation

76 KrakenParser is implemented in Python (available via PyPI as krakenparser) with several  
77 auxiliary scripts. It leverages the original KrakenTools ([Lu et al., 2022](#)) scripts for initial data  
78 reshaping and then applies its own pure-Python processing for downstream formatting. The  
79 software follows a pipeline of six main steps, which can be executed automatically in sequence  
80 (--complete mode) or run individually as needed:

- 81 1. Convert reports to MPA format: Each Kraken2/Bracken/Metaboli report (text file with  
82 taxon lines) is converted to an "MPA" table format using KrakenTools' kreport2mpa.py  
83 script. In MPA format, each row corresponds to a read and columns correspond to  
84 taxonomic ranks, allowing easy combination of multiple samples.
- 85 2. Combine MPA files: All per-sample MPA files are merged into a single master table  
86 (samples × taxa) using KrakenTools' combine\_mpa.py. This yields a matrix of raw read  
87 counts, with entries where a taxon is absent in a sample filled with zero.
- 88 3. Deconstruct taxonomic levels: The combined data is split out by rank. KrakenParser  
89 extracts separate text files for phylum, class, order, family, genus, and species counts.  
90 During this step, it can optionally isolate certain domains; for example, using --  
91 deconstruct\_viruses will produce a file of only viral species counts, ignoring other  
92 domains. Also, the default --deconstruct excludes reads classified as human to focus  
93 on microbial content.
- 94 4. Process extracted data: Each rank-specific text file is cleaned and formatted.

95        KrakenParser removes classification prefixes (like "s\_\_" for species, "g\_\_" for genus)  
96        and replaces underscores with spaces for readability. This step ensures taxon names are  
97        human-friendly (e.g. "s\_\_Escherichia\_coli" becomes "Escherichia coli").  
98        5. Convert to CSV: The cleaned text tables are converted to CSV files (comma-separated  
99        values). In this transpose operation, taxa become columns and sample identifiers become  
100      rows, yielding a standard matrix format. This structured CSV is easy to import into  
101      statistical software, spreadsheets, or R/Python data frames for further analysis.  
102      6. Calculate relative abundances: For each count table, KrakenParser can create a  
103        corresponding relative abundance table (--relabund option) by computing per-  
104        centages of total reads per sample, using the formula: Relative Abundance =  
105        
$$\left( \frac{\text{Number of individuals of taxa}}{\text{Total number of individuals of all taxa}} \right) \times 100$$
. Users can specify a threshold to group  
106        low-abundance taxa into an "Other" category. This results in a normalized profile for  
107        each sample, often more interpretable in comparative studies than raw counts.  
108      Each of these steps is exposed as a sub-command in the CLI, so advanced users can integrate  
109        KrakenParser into custom workflows. By default, running KrakenParser --complete -i  
110        <reports\_dir>/kreports executes all steps sequentially, writing outputs to a structured  
111        directory tree (with subfolders for each step). The outputs include one CSV file per rank  
112        (e.g. counts\_phylum.csv, counts\_species.csv) containing absolute read counts, and similarly  
113        named files under a csv\_relabund/ directory for percentages if requested. KrakenParser is  
114        optimized for speed and memory efficiency given the nature of the task: it processes text files  
115        line by line and uses pandas data frames for merging and calculations, which easily handle  
116        dozens of samples and tens of thousands of taxa on a standard workstation. The reliance  
117        on KrakenTools for the initial conversion ensures that the parsing logic benefits from the  
118        robustness of well-tested scripts, while the unified interface adds convenience. The tool also  
119        includes built-in help for each subcommand (-h), guiding users on required inputs and options.  
120        KrakenParser's design reflects practical needs observed in the metagenomics community -  
121        it was tested during the [2025 "Bioinformatics Bootcamp"](#) hackathon organized by ITMO  
122        University, where teams analyzing metagenomic datasets were able to obtain meaningful results  
123        in a short time thanks to KrakenParser's streamlined processing pipeline. By combining  
124        established methods with new automation, KrakenParser provides an efficient, reproducible,  
125        and user-friendly means to handle the otherwise tedious steps of post-classification data  
126        processing.  
127      KrakenParser also offers a suite of Python-based visualization tools to facilitate the interpre-  
128        tation of taxonomic profiles:  
129        ▪ Stacked Bar Plots: Utilizing matplotlib ([Hunter, 2007](#)) and pandas ([team, 2020](#)),  
130        KrakenParser can generate stacked bar plots that display the relative abundances of  
131        taxa across multiple samples. These plots provide a clear comparison of taxonomic  
132        compositions between samples.  
133        ▪ Streamgraphs: For a more dynamic representation, KrakenParser can create stream-  
134        graphs using matplotlib's ([Hunter, 2007](#)) stackplot function with a symmetric baseline.  
135        This visualization emphasizes changes in taxa abundances over a series of samples,  
136        highlighting temporal or sequential patterns.  
137        ▪ Combined Visualizations: To offer both detailed and overarching views, KrakenParser  
138        supports combined plots that integrate stacked bar plots and streamgraphs. This dual  
139        representation aids in comprehensive data analysis.  
140        ▪ Clustermaps: Employing seaborn ([Waskom, 2021](#)), KrakenParser can produce clus-  
141        termaps that perform hierarchical clustering on taxa and samples. These heatmaps reveal  
142        patterns and groupings in the data, facilitating the identification of similar taxonomic  
143        profiles.  
144      These visualization tools are accessible through the KrakenParser Python API, allowing users  
145        to customize and integrate them into their analysis workflows seamlessly.

## <sup>146</sup> Acknowledgements

<sup>147</sup> The development of KrakenParser was supported by the Russian Science Foundation (Project  
<sup>148</sup> no. 25-24-00351).

## <sup>149</sup> References

- <sup>150</sup> Breitwieser, F. P., Baker, D. N., & Salzberg, S. L. (2018). KrakenUniq: Confident and  
<sup>151</sup> fast metagenomics classification using unique k-mer counts. *Genome Biology*, *19*, 198.  
<sup>152</sup> <https://doi.org/10.1186/s13059-018-1568-0>
- <sup>153</sup> Breitwieser, F. P., & Salzberg, S. L. (2020). Pavian: Interactive analysis of metagenomics  
<sup>154</sup> data for microbiome studies and pathogen identification. *Bioinformatics*, *36*(4), 1303–1304.  
<sup>155</sup> <https://doi.org/10.1093/bioinformatics/btz715>
- <sup>156</sup> Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing in Science &*  
<sup>157</sup> *Engineering*, *9*(3), 90–95. <https://doi.org/10.1109/MCSE.2007.55>
- <sup>158</sup> Kim, D., Song, L., Breitwieser, F. P., & Salzberg, S. L. (2016). Centrifuge: Rapid and  
<sup>159</sup> sensitive classification of metagenomic sequences. *Genome Research*, *26*(12), 1721–1729.  
<sup>160</sup> <https://doi.org/10.1101/gr.210641.116>
- <sup>161</sup> Kim, J., & Steinegger, M. (2024). Metabuli: Sensitive and specific metagenomic classification  
<sup>162</sup> via joint analysis of amino acid and DNA. *Nature Methods*, *21*(6), 971–973. <https://doi.org/10.1038/s41592-024-02273-y>
- <sup>164</sup> Lu, J., Breitwieser, F. P., Thielen, P., & Salzberg, S. L. (2017). Bracken: Estimating  
<sup>165</sup> species abundance in metagenomics data. *PeerJ Computer Science*, *3*, e104. <https://doi.org/10.7717/peerj-cs.104>
- <sup>167</sup> Lu, J., Rincon, N., Wood, D. E., Breitwieser, F. P., Pockrandt, C., Langmead, B., Salzberg, S.  
<sup>168</sup> L., & Steinegger, M. (2022). Metagenome analysis using the kraken software suite. *Nature  
Protocols*, *17*, 2815–2839. <https://doi.org/10.1038/s41596-022-00738-y>
- <sup>170</sup> Menzel, P., Ng, K. L., & Krogh, A. (2016). Fast and sensitive taxonomic classification for  
<sup>171</sup> metagenomics with kaiju. *Nature Communications*, *7*, 11257. <https://doi.org/10.1038/ncomms11257>
- <sup>173</sup> Ounit, R., Wanamaker, S., Close, T. J., & Lonardi, S. (2015). CLARK: Fast and accurate  
<sup>174</sup> classification of metagenomic and genomic sequences using discriminative k-mers. *BMC  
Genomics*, *16*, 236. <https://doi.org/10.1186/s12864-015-1419-2>
- <sup>176</sup> Sczyrba, A., Hofmann, P., Belmann, P., Koslicki, D., Janssen, S., Dröge, J., Gregor, I., Majda,  
<sup>177</sup> S., Fiedler, J., Dahms, E., & others. (2017). Critical assessment of metagenome inter-  
<sup>178</sup> pretation: A benchmark of metagenomics software. *Nature Methods*, *14*(11), 1063–1071.  
<sup>179</sup> <https://doi.org/10.1038/nmeth.4458>
- <sup>180</sup> team, T. pandas development. (2020). *Pandas-dev/pandas: pandas* (latest). Zenodo.  
<sup>181</sup> <https://doi.org/10.5281/zenodo.3509134>
- <sup>182</sup> Waskom, M. L. (2021). Seaborn: Statistical data visualization. *Journal of Open Source  
Software*, *6*(60), 3021. <https://doi.org/10.21105/joss.03021>
- <sup>184</sup> Wood, D. E., Lu, J., & Langmead, B. (2019). Improved metagenomic analysis with kraken 2.  
<sup>185</sup> *Genome Biology*, *20*, 257. <https://doi.org/10.1186/s13059-019-1891-0>
- <sup>186</sup> Wood, D. E., & Salzberg, S. L. (2014). Kraken: Ultrafast metagenomic sequence classifi-  
<sup>187</sup> cation using exact alignments. *Genome Biology*, *15*(3), R46. <https://doi.org/10.1186/gb-2014-15-3-r46>