

Metabuli: sensitive and specific metagenomic classification via joint analysis of amino-acid and DNA

Martin Steinegger

martin.steinegger@snu.ac.kr

Seoul National University <https://orcid.org/0000-0001-8781-9753>

JAEBEOM KIM

Brief Communication

Keywords:

Posted Date: July 6th, 2023

DOI: <https://doi.org/10.21203/rs.3.rs-3061195/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: There is **NO** Competing Interest.

Version of Record: A version of this preprint was published at Nature Methods on May 20th, 2024. See the published version at <https://doi.org/10.1038/s41592-024-02273-y>.

Metabuli: sensitive and specific metagenomic classification via joint analysis of amino-acid and DNA

Jaebeom Kim¹ and Martin Steinegger^{1,2,3,4,✉}

¹Interdisciplinary Program in Bioinformatics, Seoul National University, Seoul, Republic of Korea

²School of Biological Sciences, Seoul National University, Seoul, Republic of Korea

³Institute of Molecular Biology and Genetics, Seoul National University, Seoul, Republic of Korea

⁴Artificial Intelligence Institute, Seoul National University, Seoul, Republic of Korea

1 **Current metagenomic classifiers analyze either DNA or**
2 **amino-acid (AA) sequences. DNA-based methods have**
3 **better specificity in distinguishing well-studied clades, but**
4 **they have limited sensitivity in detecting under-studied**
5 **clades. AA-based methods suffer the opposite problem.**
6 **To tackle this trade-off, we developed Metabuli for a joint**
7 **analysis of DNA and AA using a novel k-mer, *metamer*. In**
8 **benchmarks, Metabuli was simultaneously as specific as**
9 **DNA-based methods and as sensitive as AA-based meth-**
10 **ods. In the CAMI2 plant-associated dataset, Metabuli**
11 **covers 99% and 98% of classifications of state-of-**
12 **the-art DNA-based and AA-based classifiers, respectively.**
13 **Metabuli is available as free and open-source software for**
14 **Linux and macOS at metabuli.steineggerlab.com.**

15 **Correspondence: martin.steinegger@snu.ac.kr**

16 Metagenomics allows studying microbial communities by
17 analyzing DNA or RNA sequences directly taken from vari-
18 ous environments. Some studies aim to reveal evolutionary
19 distant organisms (e.g., in the soil (1), ocean (2) and hy-
20 drothermal vent sites (3)). Others, in the clinical field, focus
21 on detecting pathogens and emerging strains in samples from
22 patients (4), public spaces (5), and wastewater (6).

23 Identifying the origin of metagenomic reads is performed
24 by searching for similar regions in reference sequences. One
25 way to detect the similarity is to calculate local alignments
26 between the read and the reference as in MMseqs2 Taxonomy
27 (7) and MEGAN CE (8). Alternatively, alignment-free meth-
28 ods were introduced for faster classification. For instance, k-
29 mer-based tools extract fixed-length k-mers from queries and
30 references and matches them. Another type, FM-index-based
31 tools utilize the Burrows-Wheeler transformation of the refer-
32 ences to query (9, 10) k-mer matches of flexible length.

33 Metagenomic classifier needs two contrasting capabilities:
34 1) specificity for high-resolution classification of well-
35 studied clades and 2) sensitivity to detect under-studied
36 species based on known relatives in a database.

37 However, current tools suffer an inherent trade-off prob-
38 lem between specificity and sensitivity depending on the se-
39 quence type they utilize: DNA or amino-acids (AAs) (11–
40 13). DNA-based tools have better specificity as they exploit
41 point mutations to differentiate strains. AA-based tools lever-
42 age the higher conservation of AA sequences for better sensi-
43 tivity to detect homology between novel organisms and their
44 relatives in the reference, although it limits resolving close
45 taxa.

46 As a partial countermeasure, classifiers that are particu-
47 larly well-suited to the research context need to be selected

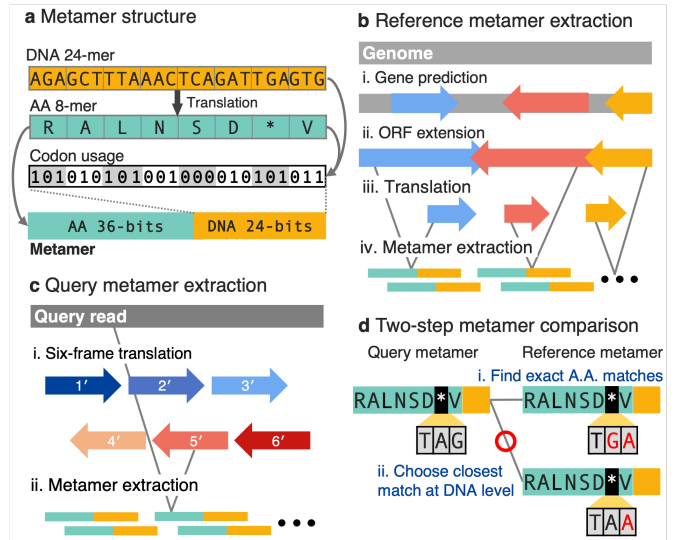


Fig. 1. Metabuli's workflow. a) A DNA fragment of 24 nucleotides is translated into eight AAs, which are encoded as an integral value encoded within 36 bits. Each AA has 1-6 synonymous codons, thus requiring three bits to store which one is seen in the fragment. b) Metabuli predicts ORFs in a genome using Prodigal and extends them to cover intergenic regions. The extended ORFs are used to extract reference metamers. c) Metabuli scans each read in six translational frames to extract query metamers. d) The metamers are compared first to find exact AA matches and subsequently to choose the closest one at the DNA level.

(11–13). However, metagenomic samples are a mixture of well- and under-studied taxa, the specificity-sensitivity trade-off inevitably restricts full sample characterization.

To address this trade-off problem, we introduce Metabuli, a method that jointly analyzes DNA sequences and their AA translation to achieve both specificity and sensitivity simultaneously (Fig. 1 and Supp. Fig. 1). In benchmarks comprising simulated reads, Critical Assessment of Metagenome Interpretation 2 (CAMI2 (15)) datasets, as well as real-world metagenomes, Metabuli consistently demonstrated top performance while DNA- and AA-based tools had fluctuating performance depending on the distance between the queried organisms and available references in the database.

To enable the joint analysis of DNA and AA sequences, Metabuli utilizes a novel k-mer structure, *metamer*, encoding a 24 nucleotide-long fragment (eight codons) in 60 bits. Its translation to AAs is encoded by 36 bits, and its codons - by 24 bits. Since an AA is coded by at most six codons, three bits per AA suffice to indicate which one is seen. This joint-encoding is more efficient compared to individual encoding, requiring only 2/3 of the bits.

During database creation, Metabuli predicts open reading frames (ORFs) using Prodigal (16). Each ORFs is extended

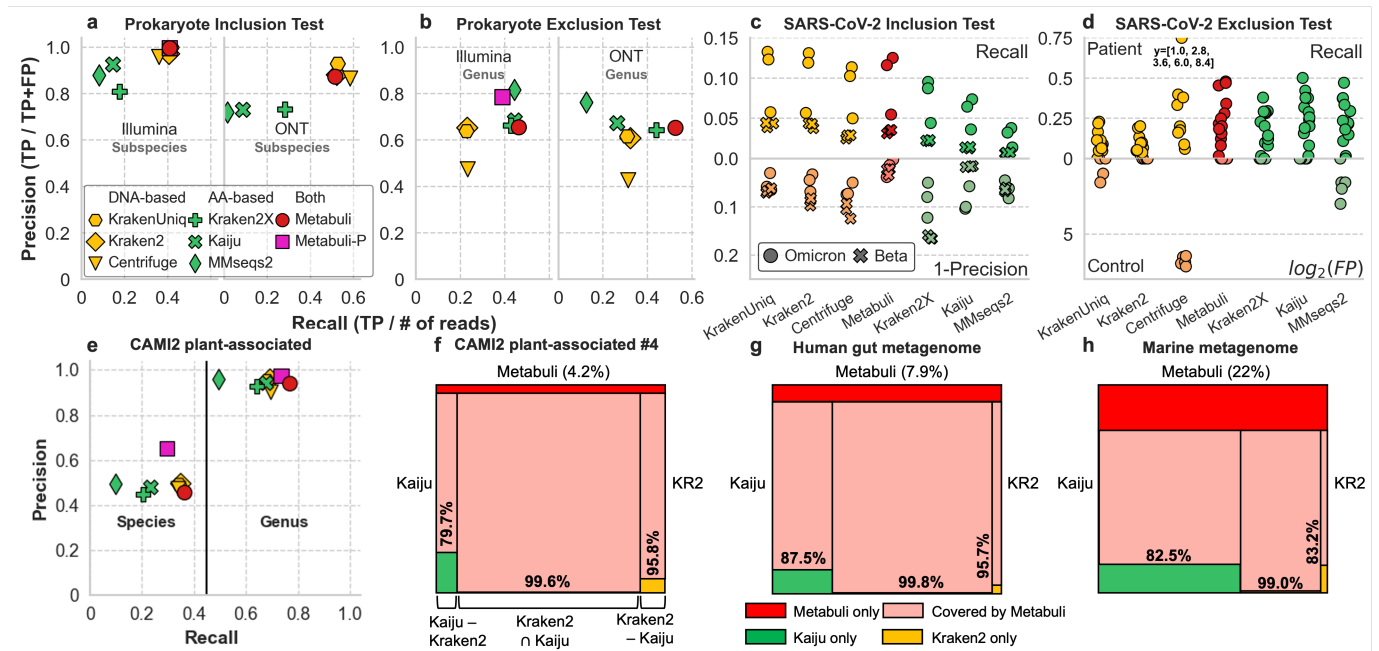


Fig. 2. Benchmark results. **a-b) GTDB benchmarks** GTDB genomes and taxonomy were used. Simulated short (Illumina) and long (ONT) reads were used. **a)** Reads were simulated from genomes present in databases. **b)** Not the queried species but their sibling species were contained in databases. **c-d) Pathogen detection tests.** RNA-seq reads from COVID-19 patients were classified. **c)** The reference included five SARS-CoV-2 variants and Viral RefSeq (14), and the reads from patients infected with either omicron or beta variant were queried. Classifications to correct and incorrect variants were counted as TP and FP, respectively. **d)** Viral RefSeq, excluding SARS-CoV-2, served as the reference. RNA-seq reads from controls (bottom) and patients (top) were classified. Classifications to Sarbecovirus, the LCA of SARS-CoV-1 and 2, were counted as TP for patient samples and as FP for controls. Centrifuge classified more reads as SARS-CoV-2 than the estimated total number of SARS-CoV-2 reads. Recall values for such cases were denoted. **e-f) CAMI2 plant-associated dataset.** **e)** GTDB genomes and the CAMI2-provided taxonomy were used for database construction. Classifications for the CAMI2-provided plant-associated reads were evaluated. **f)** The relationship among TP (at genus level) sets of Kaiju, Kraken2, and Metabuli for one of the plant-associated samples. **g-h) Real metagenomes** Human gut (g) and marine (h) metagenomic reads were classified using databases of (a). The proportion of Metabuli-only area in the union of the three tools is denoted in parentheses. **f-h)** The area is proportional to the number of reads within each panel.

71 to cover intergenic regions to cover the whole genome, these 99
 72 regions are often missed in methods utilizing only coding se- 100
 73 quences (12, 13). Notably, Metabuli only stores metamers up 101
 74 to one-third of the length of contigs. This is in contrast to AA 102
 75 classifier kAsA (17), which involves storing all k-mers from 103
 76 six frames of the entire genomes, leading to a sixfold increase 104
 77 in size. In addition, Metabuli's reference metamer list is also 105
 78 shortened by removing metamers that are redundant within 106
 79 each species. 107

80 To classify each read, Metabuli computes query metamers 108
 81 from each read and its six-frame translations, which are car- 109
 82 ried through stop codons. These are compared to reference 110
 83 metamers to find perfect AA matches for sensitivity; among 111
 84 them, matches of the lowest DNA Hamming distance are 112
 85 selected for specificity. Metabuli can quickly calculate the 113
 86 distance with a pre-computed distance matrix designed for 114
 87 metamers. The selected matches are analyzed to score can- 115
 88 didate taxa and to classify (Supp. Fig. 2). In this process, 116
 89 Metabuli-P (precision mode) uses score thresholds to reduce 117
 90 false positive and over-confident classifications (Methods, 118
 91 Supp. Fig. 3). 119

92 To compare the performance of Metabuli to state-of-the- 120
 93 art classifiers, we conducted inclusion and exclusion tests us- 121
 94 ing prokaryotes and viruses (Fig. 2a-d). In inclusion tests, 122
 95 we evaluated specificity, i.e., how well a classifier can distin- 123
 96 guish between reads from closely related organisms at lower 124
 97 taxonomic ranks. Thus, query (sub)species were present in 125
 98 the reference as well as their siblings. In contrast, exclusion 126

tests evaluated sensitivity, i.e., the ability to classify reads 99
 from a novel (sub)species based on sequences of its siblings, 100
 so the query (sub)species was removed from the reference. 101

Depending on the purpose of each test, we measured the 102
 precision (P) and recall (R) at different ranks. In inclusion 103
 tests, we measured them at the (sub)species rank, and in exclu- 104
 sion tests - at the rank of the lowest common ancestor 105
 (LCA) of each query and its siblings. When measuring at 106
 a certain rank, unclassified reads as well as reads classified 107
 at higher ranks were considered false negatives (FNs) to pe- 108
 nalize less informative classifications. Meanwhile, classifica- 109
 tions at lower ranks climbed up the taxonomy to the rank of 110
 measurement. Afterward, classifications to the correct or 111
 to the wrong taxon were counted as true positives (TPs) and 112
 false positives (FPs), respectively. 113

First, we designed a short read benchmark using the 114
 Genome Taxonomy Database (GTDB) (15). In the inclusion 115
 test, where reads were simulated from 1,191 species 116
 that had at least two subspecies in the database (19% of all 117
 species), DNA-based methods classified more reads to cor- 118
 rect subspecies than AA-based ones (Fig. 2a). AA-based 119
 methods classified less than 18% of the reads, about half 120
 of what DNA-based methods could, also with lower preci- 121
 sion. However, in the exclusion test, where reads were simu- 122
 lated from 367 species that were removed from the database, 123
 AA-based tools performed better (Fig. 2b). They classified 124
 about twice as many reads as DNA-based tools into the cor- 125
 rect genus with better precision ($R > 0.4$ for AA-based, $R <$ 126

0.25 for DNA-based). These results clearly demonstrate the pros and cons of DNA- or AA-based tools.

Next, we conducted similar tests using simulated long reads. Again, DNA-based tools outperformed AA-based ones in the inclusion test. In the exclusion test only Kraken2X, exceeded DNA-based ones. Kraken2X ignores frame information, while the other AA-tools are sensitive to frame-shifting indel errors that are more frequent in long reads (18).

Remarkably, only Metabuli achieved top-level performance in all the inclusion and exclusion tests using short and long reads. In the inclusion test, Metabuli performed as well as all DNA-based methods and outperformed all AA-based tools (Fig. 2a). Its performance was more similar to that of DNA-based tools in species rank (Supp. Fig. 4). Moreover, in the exclusion tests (Fig. 2b), Metabuli achieved the best recall with competent precision with both short and long reads. Since Metabuli scores candidate taxa using matches from multiple frames like Kraken2X, it could be robust to the frequent indels of long reads. Metabuli-P was tested only with short reads for which it is optimized, and it was the second most precise tool with comparable R to AA-based tools in the short read exclusion test.

Next, the classifiers were evaluated using real SARS-CoV-2 data for two main pathogen detection tasks: strain identification and emerging pathogen discovery, both were performed in inclusion and exclusion tests (Fig 2c-d). In the inclusion test, RNA-seq reads from six COVID-19 patients were examined to identify the culprit variant when its genome was present in databases. In contrast, only SARS-CoV-1, but not 2, was provided in the reference databases of the exclusion test.

DNA-based tools classified more reads to the culprits than the AA-based tools in the inclusion test. In the exclusion test, however, the best-performing DNA-based tool, KrakenUniq (19) missed two patient samples and made FP hits in three controls. On the other hand, AA-based tools outperformed DNA-based tools in the exclusion test, detecting up to twice as much SARS-CoV-2. However, they were worse at deciphering the specific culprit strain in the inclusion test.

Here as well, it was only Metabuli and Metabuli-P that showed robust performance in both tests. Their performance was similar, so only Metabuli is depicted in Fig. 2c-d. In the inclusion test, Metabuli classified a comparable number of reads to the culprits as DNA-based methods, even outperforming Centrifuge (10). Moreover, it achieved the best precision, classifying fewer reads to incorrect variants. In the exclusion test, it detected as many SARS-CoV-2 reads as the AA-based Kaiju (9) without any FP hits in the controls.

Next, we sought to challenge the classifiers to identify reads from datasets that contained organisms varying in their query-to-database distances, as would be the case in many real-world studies. To that end, we used query datasets from CAMI2: strain-madness, marine, and plant-associated, which have different query-to-database distances.

On the strain-madness data (Supp. Fig. 5a), Metabuli and DNA-based tools performed better than AA-based tools.

In the marine benchmark (Supp. Fig. 5b), which contains reads with larger query-to-database distances, the gap in recall became smaller and all tools showed similar precision (>0.93). For the plant-associated data with the largest query-to-database distance, tools of both types showed similar performance while Metabuli had the best sensitivity. To investigate this result, we analyzed Metabuli with respect to the genus-level TP sets of the best-performing AA- and DNA-based tools, Kaiju and Kraken2. We found that Metabuli covered 99.5% of their intersection, 76.6% of Kaiju–Kraken2, and 94.1% of Kraken2–Kaiju, which implies that Metabuli successfully joins DNA- and AA-based classifications. Moreover, about 4.2% of the total reads were correctly classified only by Metabuli. Across the three CAMI2 datasets, Metabuli-P progressively improved in precision with the growing diversity of data, with the largest improvement on the plant-associated data set (Supp. Fig. 5).

Next, we compared Kraken2, Kaiju, and Metabuli using real metagenomic data from well-studied (human gut) and under-studied (marine) environments. As real reads have no ground-truth labels, we compared the proportion of reads classified by each tool. For the human gut data (Fig. 2g), Kraken2 and Kaiju respectively classified 50% and 65% of the total. However, their classified proportion dropped significantly to 30% and 12% as query-to-database distance increased in the marine data set (Fig. 2h). On both data sets, Metabuli could classify the largest number of reads, covering 83-88% of Kaiju–Kraken2, 83-96% of Kraken2–Kaiju, and $>99\%$ of $\text{Kaiju} \cap \text{Kraken2}$.

Finally, we compared the speed, RAM usage, and database size in the prokaryote benchmarks (Supp. Table 1). All tools took less than ten minutes except for MMseqs2 Taxonomy, which spent >100 minutes. Of all, Kraken2X was the fastest and used the least RAM, also having the smallest database. Notably, because Metabuli is designed to utilize a user-specified size of RAM, it can classify reads against any size database as long as it fits in the machine's hard disk. We demonstrated this feature by measuring performance under various configurations. Metabuli was even able to complete the tasks on a notebook with just 8 GiB RAM and 8 threads (Supp. Table 1). Even though Metabuli stores both DNA and amino acid sequences, its database size was about 1.5 times that of Kraken2's probabilistic database.

In summary, Metabuli achieves high specificity and high sensitivity simultaneously by utilizing metamers to jointly analyze sequences at both DNA and AA levels. In benchmarks, only Metabuli showed robust state-of-the-art performance, while other tools sacrificed either sensitivity or specificity depending on their type and the benchmark scenario. The results demonstrate the transformative potential of Metabuli for diverse research contexts. Metabuli allows specific classifications for reads from well-studied species while not losing sensitivity for under-studied organisms. At last, Metabuli is open-source software, and ready-to-use binaries and pre-computed databases were provided (Supp. Table 2).

239 References

- 240 1. Daniel, R. The metagenomics of soil. *Nature Reviews Microbiology* **3**, 470–478 (2005).
- 241 2. Paoli, L. *et al.* Biosynthetic potential of the global ocean microbiome. *Nature* **607**, 111–118 (2022).
- 242 3. Anderson, R. E. *et al.* Genomic variation in microbial populations inhabiting the marine seafloor
- 243 at deep-sea hydrothermal vents. *Nature Communications* **8**, 1114 (2017).
- 244 4. Wilson, M. R. *et al.* Clinical metagenomic sequencing for diagnosis of meningitis and encephalitis.
- 245 *New England Journal of Medicine* **380**, 2327–2340 (2019).
- 246 5. Danko, D. *et al.* A global metagenomic map of urban microbiomes and antimicrobial resistance. *Cell*
- 247 **184**, 3376–3393 (2021).
- 248 6. Hendriksen, R. S. *et al.* Global monitoring of antimicrobial resistance based on metagenomics
- 249 analyses of urban sewage. *Nature Communications* **10**, 1124 (2019).
- 250 7. Mirdita, M., Steinegger, M., Breitwieser, F., Söding, J. & Levy Karin, E. Fast and sensitive taxonomic
- 251 assignment to metagenomic contigs. *Bioinformatics* **37**, 3029–3031 (2021).
- 252 8. Huson, D. H. *et al.* MEGAN community edition-interactive exploration and analysis of large-scale
- 253 microbiome sequencing data. *PLoS Computational Biology* **12**, e1004957 (2016).
- 254 9. Menzel, P., Ng, K. L. & Krogh, A. Fast and sensitive taxonomic classification for metagenomics with
- 255 Kaiju. *Nature Communications* **7**, 11257 (2016).
- 256 10. Kim, D., Song, L., Breitwieser, F. P. & Salzberg, S. L. Centrifuge: rapid and sensitive classification of
- 257 metagenomic sequences. *Genome Research* **26**, 1721–1729 (2016).
- 258 11. Breitwieser, F. P., Lu, J. & Salzberg, S. L. A review of methods and databases for metagenomic
- 259 classification and assembly. *Briefings in Bioinformatics* **20**, 1125–1136 (2019).
- 260 12. Simon, H. Y., Siddle, K. J., Park, D. J. & Sabeti, P. C. Benchmarking metagenomics tools for taxo-
- 261 nomic classification. *Cell* **178**, 779–794 (2019).
- 262 13. Nooij, S., Schmitz, D., Vennema, H., Kroneman, A. & Koopmans, M. P. Overview of virus metage-
- 263 nomic classification methods and their biological applications. *Frontiers in Microbiology* **9**, 749
- 264 (2018).
- 265 14. O’Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: current status, taxonomic
- 266 expansion, and functional annotation. *Nucleic Acids Research* **44**, D733–D745 (2016).
- 267 15. Meyer, F. *et al.* Critical assessment of metagenome interpretation: the second round of challenges.
- 268 *Nature Methods* **19**, 429–440 (2022).
- 269 16. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site identification.
- 270 *BMC Bioinformatics* **11**, 119 (2010).
- 271 17. Weging, S., Gogol-Döring, A. & Grosse, I. Taxonomic analysis of metagenomic data with kasa.
- 272 *Nucleic Acids Research* **49**, e68–e68 (2021).
- 273 18. Watson, M. & Warr, A. Errors in long-read assemblies can critically affect protein prediction. *Nature*
- 274 *Biotechnology* **37**, 124–126 (2019).
- 275 19. Breitwieser, F. P., Baker, D. N. & Salzberg, S. L. Krakenuniq: confident and fast metagenomics
- 276 classification using unique k-mer counts. *Genome Biology* **19**, 198 (2018).
- 277 20. Holtgrewe, M. *Mason: a read simulator for second generation sequencing data*. Dissertation, Freie
- 278 Universität Berlin, Germany (2010).
- 279 21. Ono, Y., Hamada, M. & Asai, K. Pbsim3: a simulator for all types of pacbio and ont long reads. *NAR*
- 280 *Genomics and Bioinformatics* **4**, lqac092 (2022).
- 281 22. de la Cuesta-Zuluaga, J., Ley, R. E. & Youngblut, N. D. Struo: a pipeline for building custom
- 282 databases for common metagenome profilers. *Bioinformatics* **36**, 2314–2315 (2020).
- 283 23. Youngblut, N. & Shen, W. nick-youngblut/gtdb_to_taxdump: Zenodo release (2020). URL <https://doi.org/10.5281/zenodo.3696964>.
- 284 24. Rahaman, M. M. *et al.* Genomic characterization of the dominating Beta, V2 variant carrying
- 285 vaccinated (Oxford-AstraZeneca) and nonvaccinated COVID-19 patient samples in Bangladesh: A
- 286 metagenomics and whole-genome approach. *Journal of Medical Virology* **94**, 1670–1688 (2022).
- 287 25. Lentini, A., Pereira, A., Winqvist, O. & Reinius, B. Monitoring of the SARS-CoV-2 Omicron BA.1/BA.2
- 288 lineage transition in the Swedish population reveals increased viral RNA levels in BA.2 cases. *Med*
- 289 **3**, 636–643.e4 (2022).
- 290 26. Desai, N. *et al.* Temporal and spatial heterogeneity of host response to SARS-CoV-2 pulmonary
- 291 infection. *Nature Communications* **11**, 6319 (2020).
- 292 27. Barnes, S. J. *et al.* Metagenome-assembled genomes from photo-oxidized and nonoxidized oil-
- 293 degrading marine microcosms. *Microbiology Resource Announcements* e00210–23 (2023).
- 294 28. Lu, J. *et al.* Metagenome analysis using the kraken software suite. *Nature Protocols* **17**, 2815–2839
- 295 (2022).
- 296

297 Acknowledgements

298 The authors wish to thank Eli Levy Karin from ELKMO for
299 the scientific feedback and careful reading of the manuscript;
300 Johannes Söding for discussions on metamer encoding; Milot
301 Mirdita for the usability improvement of the software; Sebas-
302 tian Jaenicke for voluntary examination of the software; and
303 Minjae Kim for feedback on the manuscript.

304 M.S. acknowledges support from the National Research
305 Foundation of Korea (grants 2019R1A6A1A10073437,
306 2020M3A9G7103933, 2021R1C1C102065, and
307 2021M3A9I4021220), the Samsung DS Research Fund, and
308 the Creative-Pioneering Researchers Program through Seoul
309 National University.

310 Author contributions

311 J.K. and M.S. designed the research, developed the software,
312 performed analysis, and wrote the manuscript.

313 Competing interests

314 The authors declare no competing interests.

315 Methods

316 Simulated read generation

317 To simulate paired-end short reads used in synthetic bench- 371
318 marks, we used the *mason_simulator* module of Mason2 372
319 (20). The reads were 150 nt in length and included simu- 373
320 lated errors at rates of 0.11% for mismatches, 0.005% for 374
321 insertions, and 0.005% for deletions. These error rates were 375
322 based on the performance of the NovaSeq 6000 sequencer. 376
323 As with Mason2's default settings, the mismatch probability 377
324 at the beginning and end of the reads was set to 0.5% and 378
325 0.22%. When provided to MMseqs2 Taxonomy, simulated 379
326 reads were concatenated with 'NN' as it does not support 380
327 paired-end reads. In the case of long reads, we used PBSIM3 381
328 (21) to simulate reads of Oxford Nanopore Technologies 382
329 with options; `--strategy wgs --method errhmm` 383
330 `--errhmm ERRHMM-ONT.model --depth 3.` 384

331 GTDB

332 The GTDB was used for several benchmarks as well as 387
333 for the calibration of Metabuli-P as it provides phylogenet- 388
334 ically consistent taxonomy based on genomic distance mea- 389
335 sures. For these, we started with a subset of GTDB R202 390
336 consisting of 258,406 genomes from 47,894 species clus- 391
337 ters. We used the *GTDB_metadata_filter.R* module in 392
338 the pipeline Struo (22) to obtain a list of 22,973 genomes 393
339 that were assembled at the level of complete genome or chro- 394
340 mosome, had CheckM completeness >90, and had CheckM 395
341 contamination <5. The filtered genomes were downloaded 396
342 using Struo's *genome_download.R* module, and 22,819 397
343 successfully downloaded genomes of 6,186 species were 398
344 used. NCBI-style taxonomy dump files for the GTDB were 399
345 generated by *gtdb_to_taxdump* (23) module. The proteome 400
346 corresponding to each genome was computed by Prodigal 401
347 with default settings.

348 Metabuli: Database creation

349 Metabuli builds a reference database of computed metamers 404
350 from nucleotide sequences following the procedure below 405
351 (Supp. Fig. 1a-e). 406

352 **ORF prediction and extension.** Metabuli utilizes Prodigal 407
353 for ORF prediction in reference sequences. To enhance the 408
354 prediction process's efficiency, we implemented three opti- 409
355 mizations. 1) Metabuli bins reference sequences by species 410
356 in separate FASTA files, then it trains Prodigal once for each 411
357 species using the longest sequence of the species' bin before 412
358 predicting genes. This approach significantly reduces train- 413
359 ing time, considering the presence of multiple assemblies for
360 a single species. 2) We narrowed down the calculation range 414
361 of Prodigal's dynamic programming during both the training 415
362 and prediction steps. While this adjustment may cause Prodi- 416
363 gal to miss very long genes, it effectively reduced runtime by 417
364 half in tests performed on an *Escherichia coli* genome. 3) 418
365 We parallelized the training and prediction processes by dis- 419
366 tributing jobs for species bins across multiple threads, further 420
367 accelerating computation. After the gene prediction, genes 421
368 that are fully nested in longer ones are removed. The ORFs 422

369 of the remaining genes are extended to cover all intergenic
370 regions while maintaining the predicted translational frame.

Reference metamer calculation and compression. Metabuli computes reference metamers from the extended ORFs and their translations. All computed metamers are sorted numerically and then by their associated species ID. Because metamers encode amino acids in the leading significant bits, metamers encoding the same amino acid sequence are placed consecutively after sorting, and within them, they are grouped by codon usage, followed by their associated species ID. Then, redundant metamers from the same species are removed, retaining only one of them (Supp. Fig. 1c). The reduced metamer list is then further compressed as follows (Supp. Fig. 1d). The full numerical value of the first metamer is stored. For all other metamers on the list, only the increment value from the previous metamer is stored. The 64-bit encoding of the first metamer and the increments are then scanned as four slices of 15 bits each (the last four bits are unused). The slice of the least significant bits and any slice where some of the bits are turned on are copied and stored in 16 bits with one extra bit for an *end* flag. The end flag indicates whether the copied slice was the last one to be saved from a specific 64-bit value (where 1 = the copied slice is the last one). The optimal case is when only one slice is stored per metamer, yielding a compression ratio of four. The more reference metamers there are, the smaller the increments between consecutive ones tend to be, so the compression rate becomes closer to four. For example, when Metabuli was used to create a database from genomes of NCBI RefSeq release 217 (~1.1TB), the compression rate was about three. Throughout this procedure, the reference sequence ID associated with each metamer is stored alongside it as well as information concerning metamer redundancy.

402 Metabuli: Database decompression and usage

403 The values of the first metamer and the increments can be 404
405 computed back from the stored compression by concatenat-
406 ing corresponding slices in a 64-bit data type. From the sec-
407 ond metamer, their values are sequentially calculated by sum-
ming up each increment.

408 Metabuli: Classification

409 **Metamer match search.** Query metamers from reads are 410
411 sorted and compared to the reference metamer list to find
412 matches (Supp. Fig. 1f-g). Because both query and refer-
413 ence metamers are sorted, a single iteration through the lists
is enough to find all matches.

Calculating Hamming distance. After a query metamer is matched with reference metamers that are identical to it on the AA level, the closest matches are selected based on their DNA Hamming distance to the query. The distance between query and reference metamers is calculated using a Hamming distance lookup table (Supp. Fig. 2a-b). In this table, the 3-bit representations of any pair of synonymous codons are used as indices to retrieve their distance. The distances of a match are summed up when the total DNA

423 Hamming distances of matches are compared to choose the 478
424 closest metamer match (Supp. Fig. 2c). 479

425 **Computing sequence similarity and assigning taxonomy.** 481
426 The matched metamers of each read are grouped by genus 482
427 and species and examined by their coordinates on the read. 483
428 For each species, only matches within a minimum of four 484
429 consecutive matches are used to reduce the risk of random 485
430 matches. Two matches are considered consecutive when 1) 486
431 their query metamers are extracted from positions that differ 487
432 by 3 nt in the same translational frame, and 2) the Hamming 488
433 distances within the overlapping region are identical. Such 489
434 matches to each genus are aligned to the query to compute the 490
435 sequence similarity score between the query and the genus. 491
436 The score is calculated based on the number of identical AAs, 492
437 the Hamming distances, and the query length (Supp. Fig. 1h 493
438 and Supp. Fig. 2c). Next, Metabuli assigns the read to the 494
439 genus of the highest sequence similarity score. If more than 495
440 one genus has scored the highest, the query is classified as 496
441 the LCA of the best-scoring genera. Similarly, the matches 497
442 found from the assigned genus are grouped by each species 498
443 to assign the query to the species of the highest sequence sim- 499
444 ilarity (Supp. Fig. 1i). 500

445 **Metabuli: Metabuli-P** 501

446 Notably, as with other short k-mer-based classifiers, rely- 502
447 ing on few matches can often lead to false positive or over- 503
448 confident classifications. False positive classification occurs 504
449 mainly when the matched region is short. The similarity be- 505
450 tween a pair of sequences is expected to be higher if the 506
451 pair belongs to the same lower taxonomic rank (rather than 507
452 a higher rank). Over-confidence occurs when a read is class- 508
453 ified at lower ranks like species or subspecies with not enough 509
454 sequence similarity. To address this, Metabuli's precision 510
455 mode (Metabuli-P) uses two sequence similarity thresholds 511
456 to avoid false and overconfident classifications. These thresh- 512
457 olds were set based on similarity score distributions within 513
458 prokaryotic and viral genera and species (Supp. Fig. 3).

459 **Distribution of sequence similarity scores.** We investigated 514
460 the distribution of sequence similarities underlying TP and 515
461 FP classifications using prokaryotes and viruses. Prokary- 516
462 otic and viral species were identified based on two crite- 517
463 ria: 1) there was at least one other species belonging to the 518
464 same genus in the database, and 2) the database contained 519
465 genomes of at least two of their subspecies. For prokary- 520
466 otes, we could find 435 species, from the 22,819 GTDB 521
467 genomes, that met the two criteria. We then designed two 522
468 settings: subspecies-exclusion and species-exclusion. In the 523
469 subspecies-exclusion, for each of the 435 species, one sub- 524
470 species was included in the reference database while one of 525
471 its sibling subspecies was excluded from it and used to sim- 526
472 ulate query reads. In the case of species-exclusion, the same 527
473 database was used, and for each of the 435 species a ran- 528
474 dom sibling species from the same genus was used to gen- 529
475 erate query reads. In both settings, 45,000 paired-end reads 530
476 for each query genome were simulated using Mason2 as de- 531
477 scribed above. In the case of viruses, we used NCBI tax-

onomy and Viral RefSeq. We could not find enough viral
species fulfilling both criteria. Therefore, for the subspecies-
exclusion setting, we applied the second criterion to find
211 species with at least two subspecies. In the case of
the species-exclusion setting, the first criterion was applied
to find 889 genera that have at least two species. In both
settings, 10,000 paired-end reads were simulated from each
query genome. Then, we used Metabuli to classify query
reads in the various test settings and examined the sequence
similarity scores underlying the TP or FP classifications.

Determining thresholds. Examination of the sequence simi-
larity distributions revealed that FP's relative frequency peaks
under sequence similarity of 0.1 (Supp. Fig. 3). Further-
more, the vast majority: 89.2-99.7% of all TPs are associated
with a sequence similarity score greater than 0.15 (Supp. Fig.
3a-d), while many FPs (29.4-59.7%) are associated with a
lower score (Supp. Fig. 3e-h). Therefore, Metabuli-P is set
to leave a query as unclassified if its best genus-level simi-
larity score is lower than 0.15. In the subspecies-exclusion
settings, 97.0% (prokaryote) and 82.2% (virus) of the TPs are
associated with a similarity score greater than 0.5 while only
14.6% (prokaryote) and 57.4% (virus) of the FPs scored as
high. Thus, Metabuli-P is set to classify a read at the species
level or a lower rank only if it has a similarity score of > 0.5
to at least one species.

Prokaryote benchmarks

Inclusion test. We examined the 22,819 complete genome or
chromosome level assemblies in the GTDB by their species
and identified 1,626 species that had at least two subspecies
with a genome in the database. Of these, 435 species were
used for the score threshold setting of Metabuli-P (Supp. Fig.
3). The remaining (1,191) contributed two subspecies each,
from which 6,150 paired-end reads were simulated with Ma-
son2 (~15M reads in total). Each of the same genomes was
also used to simulate ONT reads of 3X depth using PB-
SIM3. Performance metrics were measured at species and
subspecies ranks.

Exclusion test. The 22,819 GTDB genomes were examined
by their genera. We identified 802 genera, which had at
least two species with a genome in the database. Of these,
435 were used for the score threshold setting of Metabuli-P
(Supp. Fig. 3). The remaining 367 genera were used for the
exclusion test. In this setting, ~50,000 reads were simulated
from each species (~20M reads in total) using Mason2. PB-
SIM3 was used to simulate ONT reads of 3X depth from each
of the species. Performance was measured at genus rank.

Pathogen detection benchmarks

Inclusion test. Reference databases were built using genomes
from NCBI Viral Refseq and five SARS-CoV-2 variant
genomes (alpha, beta, delta, gamma, and omicron). We man-
ually included these variants as children of SARS-CoV-2 to
the NCBI taxonomy database. Two sets of RNA sequencing
data from COVID-19 patients were used as query reads. One

531 set was prepared from patients infected by the beta variant 584
532 (24), and the other - by the omicron variant (25). 585

533 **Exclusion test.** The database for each tool was constructed 586
534 using the taxonomy of NCBI and the genomes of Viral Ref- 587
535 Seq, excluding all SARS-CoV-2 sequences. Due to this 588
536 exclusion, SARS-CoV-1 is the closest relative in the refer- 589
537 ence database to any variant of SARS-CoV-2. RNA-seq data 590
538 from SARS-CoV-2 patients and controls prepared in a host- 591
539 response study were used as query reads (26). The estimated 592
540 number of SARS-CoV-2 reads in each sample was calculated 593
541 by multiplying the total number of RNA-seq reads by the 594
542 reads per million (RPM) of reads aligned to the SARS-CoV-2 595
543 genome. The RPM values were taken from the original study. 596

544 **CAMI 2 benchmarks**

545 We used paired-end reads of strain-madness, marine, and 597
546 plant-associated datasets and taxonomy provided in CAMI2
547 (15). In the case of CAMI2-provided reference databases for
548 DNA- and AA-based tools (*nt* and *nr*), where there are no
549 one-to-one relationships between their entries, it is possible
550 to encounter under- or over-representation of some taxa. This
551 discrepancy can lead to a potentially unfair comparison be-
552 tween the two groups of classifiers. To replace the CAMI2-
553 provided databases, we used the reference genomes and pro-
554 teomes in the prokaryote inclusion test. The references and a
555 mapping from accessions to taxonomic IDs used in CAMI2
556 were provided to each classifier for database creation. Metab-
557 uli, Centrifuge, and KrakenUniq used 7,318 genomes, which
558 together with two additional genomes were used for Kraken2,
559 Kraken2X, Kaiju, and MMseqs2. CAMI2 provides 10, 21,
560 and 100 query samples for the marine, plant-associated, and
561 strain-madness benchmarks, respectively. To reduce the run-
562 time of the benchmarks, we took all, every second, and every
563 tenth query samples, respectively. We also used the CAMI2-
564 provided ground truth labels for each read. When measuring
565 performance at the species and genus ranks, we ignored clas-
566 sifications for reads whose ground truth taxon is at a higher
567 rank than the rank of measurements.

568 **Benchmarks with real metagenomes**

569 We challenged the classifiers on two distinct metagenomes:
570 one obtained from a well-studied environment, specifically
571 a human gut sample (SRR24315757), and the other from a
572 less-studied environment, a marine sample (SRR23604821)
573 (27). The same GTDB databases as in the prokaryote inclu-
574 sion test were used.

575 **Resource measurement**

576 Maximum RAM usage (maximum resident set size) and
577 elapsed time of each tool were measured with the GNU `time`
578 `-v` command. The average performance over five repeated
579 measurements is reported (Supp. Table 1).

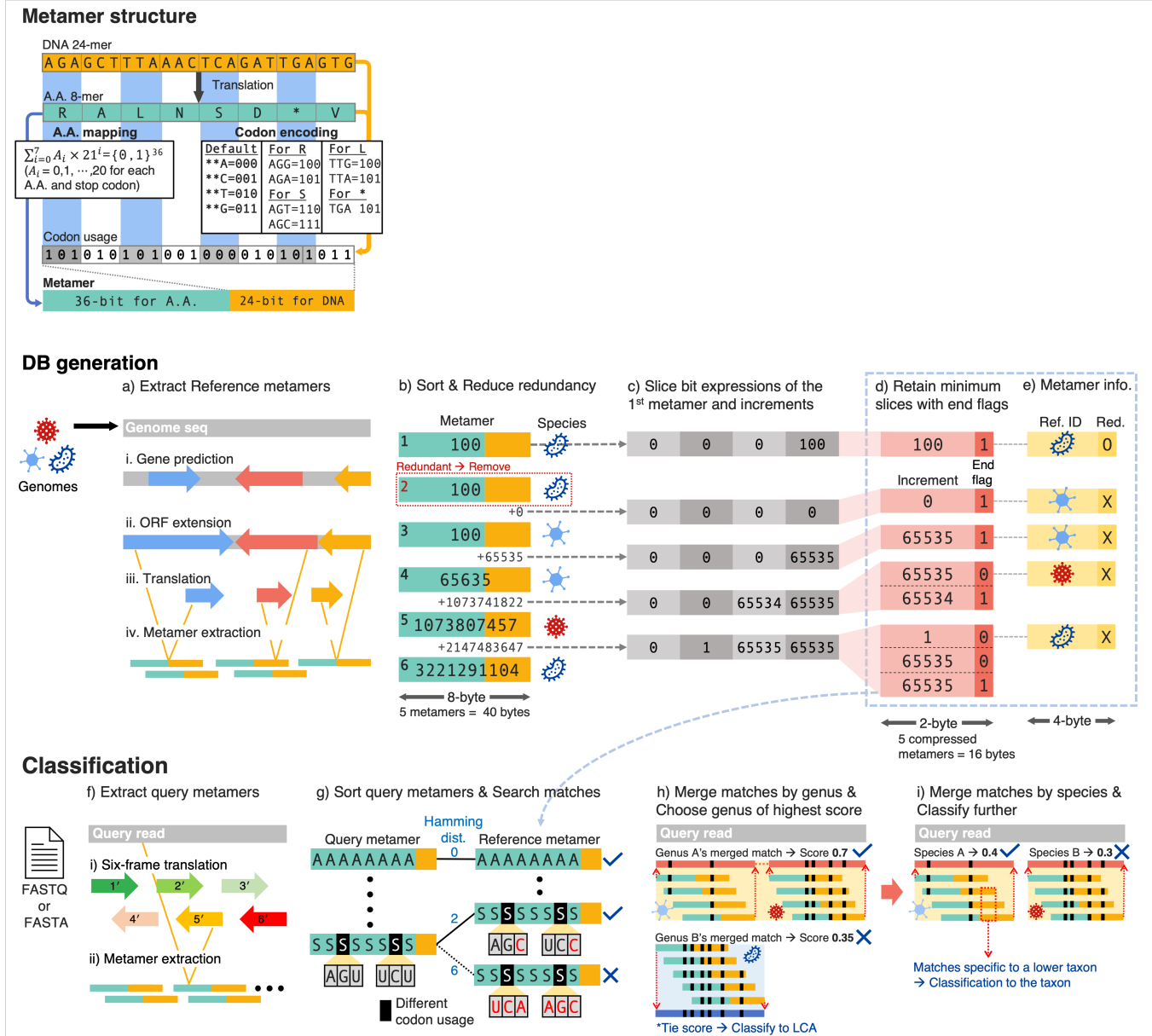
580 **Software versions and options**

581 All benchmarks were performed with Kaiju v1.9, Kraken2
582 v2.1.2, KrakenUniq v0.7.3, Centrifuge v1.0.4, and MM-
583 seqs2 v13.45111. We run Centrifuge with `-k 1` option to

report at most one classification per read. For Kraken2,
584 `--minimum-hit-groups` was set as 3 following a recom-
585 mended usage (28). Struo v0.1.7 was used to download
586 genomes and make taxonomy dump files for GTDB bench-
587 marks. Mason_simulator v2.0.9 and PBSIM3 v3.0.0 were
588 used to simulate query reads. 589

590 **Computing resource**

591 For the resource measurement, we used a server and a Mac-
592 Book Air. The server was equipped with a 64-core AMD
593 EPYC 7742 CPU and 1TB of RAM, and the MacBook Air
594 (2020) had 8GB RAM and an Apple M1 chip (8-core CPU
595 with 4 performance cores and 4 efficiency cores). A server
596 with 2×64-core AMD EPYC 7742 CPUs and 2TB of RAM
597 was used for other benchmarks.

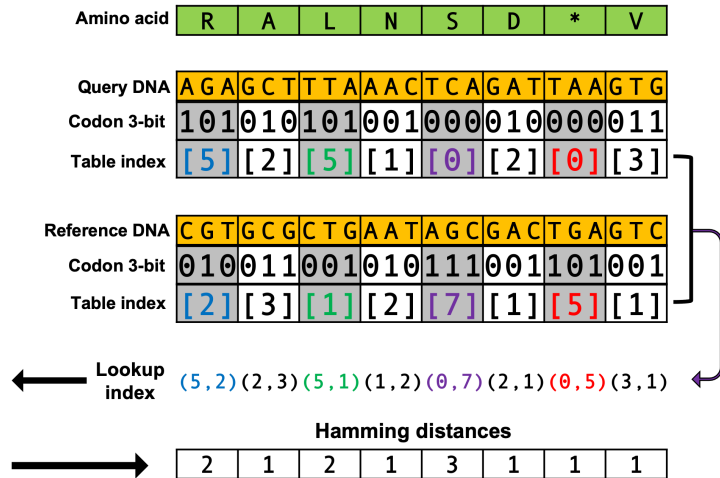


Supplementary Fig. 1. Metabuli's workflow and metamer structure. **Metamer structure.** An AA 8-mer and the codon usage of each AA are stored in a metamer using 60 bits. The codon encoding table shows how synonymous codons of each AA are mapped to 3-bit encodings. **Database generation.** a) Metabuli builds a database from genomes in FASTA format. It predicts ORFs using Prodigal and extends them to cover intergenic regions. The extended ORFs and their translations are used to compute reference metamers. b) The computed metamers are sorted numerically and redundant ones from the same species are removed. c) The 60-bit expression of the first metamer and each difference (increment) between two consecutive metamers are scanned as four 15-bit slices. d) The slice of the least significant bits is stored, followed by all other non-zero slices. An end flag is added to each slice to indicate if it was the last one to be stored from the 60-bit expression. This allows grouping the slices by the 60-bit expression they correspond to. e) The reference ID and redundancy of each metamer are stored in a separate list. **Classification.** f) Metabuli takes query files in FASTA or FASTQ format. It scans each read in six frames and computes metamers from the DNA fragments and their translations. g) Query metamers are sorted and compared to reference metamers to find perfect AA matches. Among these, matches with the smallest DNA Hamming distance are selected (Supp. Fig. 2a-b). h-i) The matches of each genus are aligned to the query to score the genus, and then matches from the best genus are grouped by species to find the best species (Supp. Fig. 2c). Matches specific to a lower rank are used for lower-rank classifications.

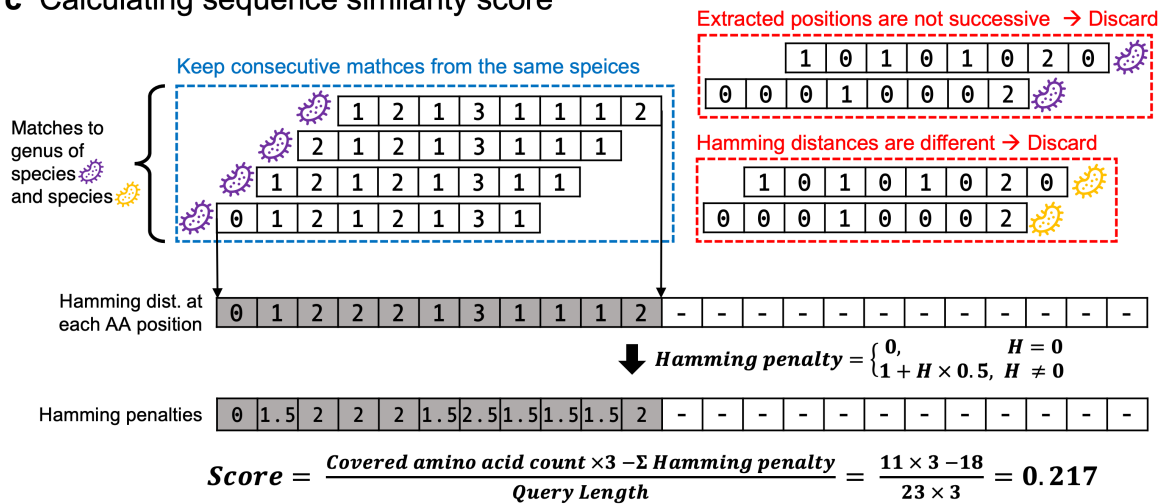
a Hamming distance matrix

	0	1	2	3	4	5	6	7	
Other	**A	**C	**T	**G	---	---	---	---	
Arg	CGA	CGC	CGT	CGG	AGG	AGA	---	---	
Leu	CTA	CTC	CTT	CTG	TTG	TTA	---	---	
Ser	TCA	TCC	TCT	TCG	---	---	AGT	AGC	
Stop	TAA	---	---	TAG	---	TGA	---	---	
0	*CCTT *GTCA AAAAA	0	1	1	1	2	1	3	3
1	*CCT- *GTC- CCCC-	0	0	1	1	2	2	3	2
2	*CCT- *GTC- TTTT-		0	0	1	2	2	2	3
3	*CCTT *GTCA GGGGG			0	0	1	2	3	3
4	-AT-- -GT-- -GG--				0	0	1	-	-
5	-AT-T -GT-G -AA-A					0	-	-	
6	---A- ---G- ---T-						0	1	
7	---A- ---G- ---C-								0

b Calculating Hamming dist. of matches

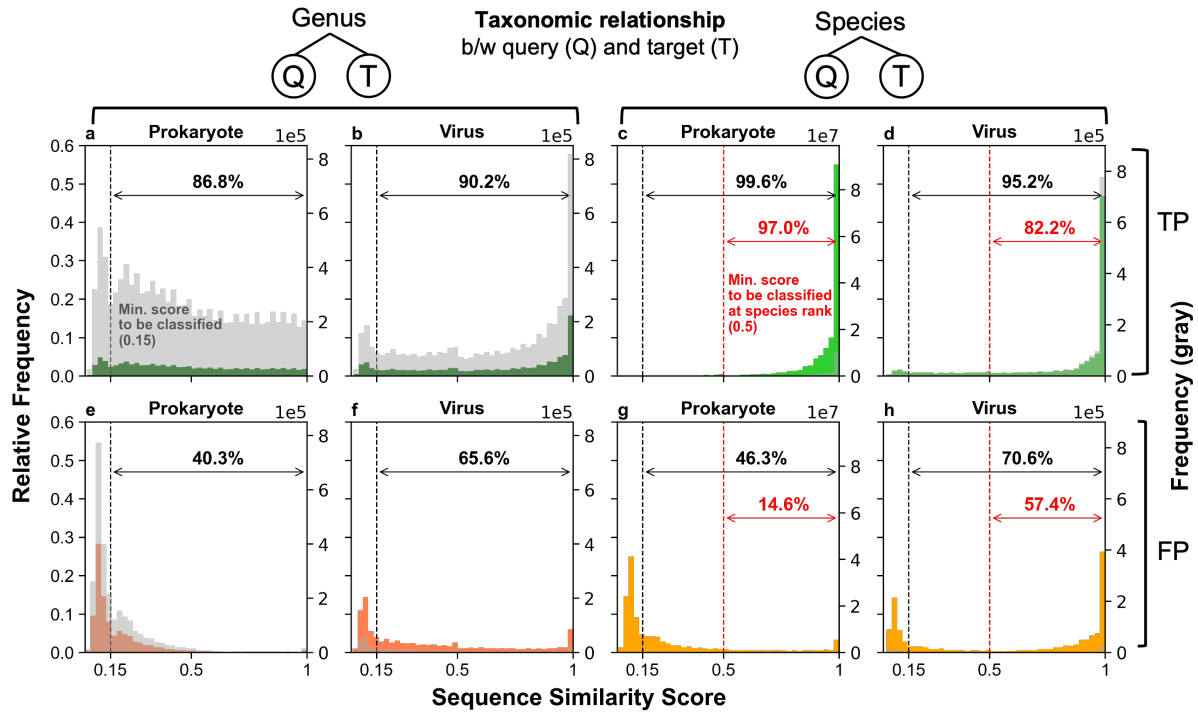


c Calculating sequence similarity score

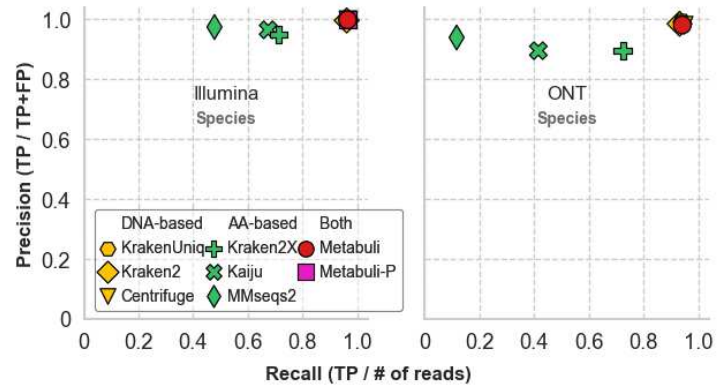


Supplementary Fig. 2. Calculating Hamming distance. a) The Hamming distance lookup table stores the distance between two codons of an identical AA pair in an 8 by 8 matrix. b) An example of Hamming distance calculation. An AA sequence (top, green) can be a translation of two different DNA sequences (DNA 1 and 2, orange). The 3-bit codon encodings for the same AA are used to index the Hamming distance lookup table. The Hamming distances of 8 codon pairs are summed up to get the total Hamming distance. c) Matches to a genus are aligned along the translated query, and Hamming distance at each position is pooled (See Supplementary Fig. 1g-h). The score of the taxon is calculated using the number of covered amino acids and Hamming penalty.

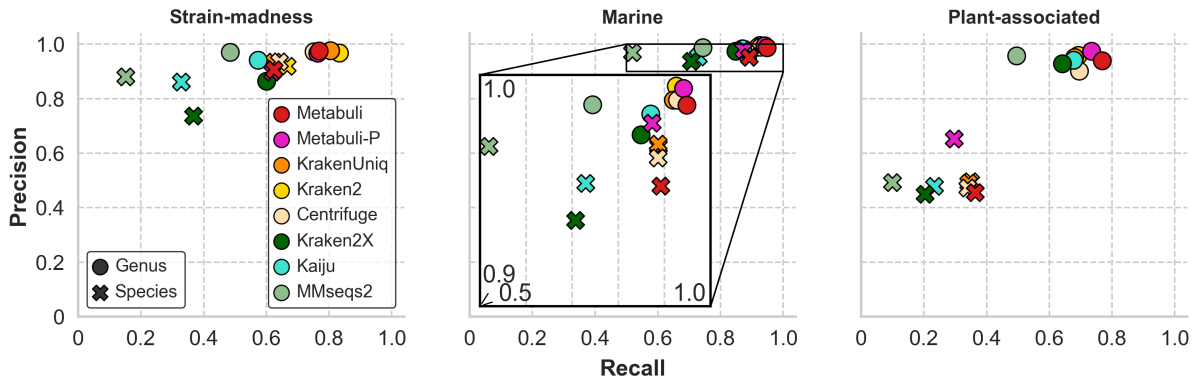
Determining Metabuli-P's score thresholds based on microbial data



Supplementary Fig. 3. Sequence similarity score distribution. The distribution of sequence similarity scores was examined in prokaryotic and viral data (full details in Methods). The thresholds for classification are marked as dashed lines. These thresholds were selected because most TP classifications were made with sequence similarity that is greater than these thresholds, while many of the FPs have lower sequence similarity. **a-h)** Setting a threshold of 0.15 as the minimal sequence similarity for classification removes 53.7-59.7% of FP prokaryotic classifications and 29.4-34.4% of the viral FPs while retaining 86.8-99.6% of all TPs. **c-d)** Out of species-level classifications, 97.0% (prokaryote) and 82.2% (virus) of TPs have sequence similarity score > 0.5. So Metabuli-P has a threshold of 0.5 as the minimal sequence similarity for species-level classification to avoid over-confident low-rank classification. A similar threshold could not be identified for the genus level (a-b).



Supplementary Fig. 4. Prokaryote inclusion test results at species rank Precision and recall of tools in the benchmarks of Fig. 2a were measured at species rank.



Supplementary Fig. 5. Benchmarks using CAMI2's strain-madness and marine dataset GTDB genomes and the CAMI2-provided taxonomy were used for reference construction. CAMI2-provided queries of strain-madness, marine, and plant-associated datasets were classified by each tool, and metrics were measured at the species and genus ranks.

Supplementary Table 1. Speed and memory usage

Software	GTDB inclusion test			GTDB exclusion test		
	DB size (GiB)	RAM (GiB)	Time (sec)	DB size (GiB)	RAM (GiB)	Time (sec)
Kraken2	43	44	24	41	42	55
KrakenUniq	309	306	169	294	292	272
Centrifuge	40	41	218	39	41	247
Kraken2X	11	12	26	10	12	47
Kaiju	39	41	145	38	41	582
MMseqs2	37	174	6075	34	173	6606
Metabuli 32GiB	69	27	525	66	22	512
Metabuli 64GiB	-	47	484	-	37	465
Metabuli 128GiB	-	91	473	-	68	448
Metabuli 256GiB	-	173	480	-	129	450
Metabuli MacBook 6GiB	-	5	5640	-	4	6680

* About 15M and 20M of 150nt paired-end reads were used in the inclusion and exclusion test

* Metabuli has an `--max-ram` option that limits maximum RAM usage. Here, runs with the option set as 6, 32, 64, 128, or 256 GiB were presented.

* All runs utilized 32 threads except for "Metabuli Macbook", which used 8 threads.

Supplementary Table 2. Pre-computed databases

Name	Size (GiB)	Data	Link
GTDB	81.2	Complete genome or chromosome level assemblies in GTDB207 and a human genome. GTDB's taxonomy was used.	https://metabuli.steineggerlab.workers.dev/gtdb207+human.tar
RefSeq	115.6	Complete Genome or Chromosome level assemblies of virus and prokaryotes in RefSeq (2023-04-04) and human genome	https://metabuli.steineggerlab.workers.dev/refseq_complete_chromosome+human.tar
RefSeq217	480.5	Refseq release 217 and human genome	https://metabuli.steineggerlab.workers.dev/refseq_release217+human.tar
RefSeq_virus	1.5	Genomes of Viral RefSeq	https://metabuli.steineggerlab.workers.dev/refseq_virus.tar

* For the human genome, GRCh38.p14 is used.

* For GTDB database, genomes with CheckM Completeness > 90 and CheckM Contamination < 5 were used.

* Taxonomy of GTDB was edited to include a human taxon.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [supplementarytables.xlsx](#)