

# UC San Diego

## UC San Diego Previously Published Works

### Title

Memory-bound k-mer selection for large and evolutionarily diverse reference libraries.

### Permalink

<https://escholarship.org/uc/item/8dg689rd>

### Journal

PCR methods and applications, 34(9)

### Authors

Şapcı, Ali

Mirarab, Siavash

### Publication Date

2024-10-11

### DOI

10.1101/gr.279339.124

Peer reviewed

# Memory-bound $k$ -mer selection for large and evolutionarily diverse reference libraries

Ali Osman Berk Şapcı<sup>1</sup> and Siavash Mirarab<sup>1,2</sup>

<sup>1</sup>Bioinformatics and Systems Biology Graduate Program, University of California, San Diego, California 92093, USA; <sup>2</sup>Department of Electrical and Computer Engineering, University of California, San Diego, California 92093, USA

Using  $k$ -mers to find sequence matches is increasingly used in many bioinformatic applications, including metagenomic sequence classification. The accuracy of these downstream applications relies on the density of the reference databases, which are rapidly growing. Although the increased density provides hope for improvements in accuracy, scalability is a concern. Reference  $k$ -mers are kept in the memory during the query time, and saving all  $k$ -mers of these ever-expanding databases is fast becoming impractical. Several strategies for subsampling have been proposed, including minimizers and finding taxon-specific  $k$ -mers. However, we contend that these strategies are inadequate, especially when reference sets are taxonomically imbalanced, as are most microbial libraries. In this paper, we explore approaches for selecting a fixed-size subset of  $k$ -mers present in an ultra-large data set to include in a library such that the classification of reads suffers the least. Our experiments demonstrate the limitations of existing approaches, especially for novel and poorly sampled groups. We propose a library construction algorithm called  $k$ -mer RANKer (KRANK) that combines several components, including a hierarchical selection strategy with adaptive size restrictions and an equitable coverage strategy. We implement KRANK in highly optimized code and combine it with the locality-sensitive hashing classifier CONSULT-II to build a taxonomic classification and profiling method. On several benchmarks, KRANK  $k$ -mer selection significantly reduces memory consumption with minimal loss in classification accuracy. We show in extensive analyses based on CAMI benchmarks that KRANK outperforms  $k$ -mer-based alternatives in terms of taxonomic profiling and comes close to the best marker-based methods in terms of accuracy.

[Supplemental material is available for this article.]

The number of genomes available in public repositories has been growing dramatically in recent years, especially due to increased sequencing of microbial and viral species. Of the 322,193 bacterial species deposited on RefSeq since 2000, >15% have been added in 2023 alone. A major benefit of having access to all these genomes is to build ultra-large reference libraries with the potential to search new *query* sequences against them. Such libraries can be used to classify reads from a metagenomic sample and to detect contaminants. It has been long appreciated that the accuracy of these downstream applications in metagenomics relies on having access to dense reference libraries as classification accuracy suffers when the query is distant from all reference genomes (Pachiadaki et al. 2019; von Meijenfeldt et al. 2019; Liang et al. 2020; Rachtman et al. 2020). Many authors have recently attempted to create ultra-large genomic reference sets (e.g., Parks et al. 2018; McDonald et al. 2024), often used with marker genes (Asnicar et al. 2020; Balaban et al. 2024). Using these data sets with  $k$ -mers-based methods, however, runs into a mundane but key limitation—the memory needed to use these genomes as reference libraries.

Using  $k$ -mers to build ultra-large libraries and classify reads has been a promising approach, as evident from the success of Kraken 2 (Wood et al. 2019). Kraken 2 and other methods such as CLARK (Ounit and Lonardi 2016) and CONSULT-II (Şapcı et al. 2024) extract  $k$ -mers for some large  $k$  (e.g., 31) from the reference set, use some strategy to choose a subset of  $k$ -mers, and build a searchable data structure (e.g., a hash table) of the selected  $k$ -mers

and their taxonomic associations. At the time of the query, the entire data structure is loaded into the memory,  $k$ -mers from reads are extracted and searched against the data structure, and a read is classified at some taxonomic rank if certain heuristic conditions (specific to each method) are met. It is easy to store millions or even a few billion  $k$ -mers in memory but modern data sets can include tens or hundreds of billions of  $k$ -mers. For example, Şapcı et al. (2024) were able to fit 8 billion 32-mers and their taxonomic labels into  $\approx 140$  GB of memory, but these had to be subsampled (arbitrarily) from 20 billion 32-mers available in the moderate-size database of 10,575 genomes produced several years ago by Zhu et al. (2019). Trade-offs of accuracy and scalability imposed by the need to subsample  $k$ -mers will be increasingly felt by developers of  $k$ -mer matching methods and those aiming to build ultra-large reference sets.

For the  $k$ -mer-based methods to benefit from the ever-growing set of available genomes, we need improved methods of subsampling  $k$ -mers in a memory-bound fashion. Of course, we are not the first to note the need to subsample  $k$ -mers. Minimization is the standard technique that selects some among adjacent  $k$ -mers (Roberts et al. 2004; Zheng et al. 2023) and is adopted by most tools. However, we argue that beyond minimization, which focuses on the redundancy of adjacent  $k$ -mers, we need to consider the evolutionary dimension. This direction can be stated as a computational problem: We are given a taxonomic tree  $\mathcal{T}$  and a set of genomes  $\mathcal{G}$  labeled by the taxonomy. Let  $\mathcal{K}$  be the set of all distinct  $k$ -mers across all genomes  $g \in \mathcal{G}$ . We seek a subset of  $\mathcal{K}$  with size  $M \leq |\mathcal{K}|$  such that the reference taxonomy is represented well

**Corresponding author:** smirarab@ucsd.edu

Article published online before print. Article, supplemental material, and publication date are at <https://www.genome.org/cgi/doi/10.1101/gr.279339.124>. Freely available online through the *Genome Research* Open Access option.

© 2024 Şapcı and Mirarab. This article, published in *Genome Research*, is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

with the selected  $k$ -mers. Users can adjust  $M$  to control the memory budget (see Methods for details).

The problem statement leaves the meaning of “well-represented” unspecified as many criteria can be considered, and the choice of the criterion forms the basis of each method. A notable attempt in this direction is pioneered by Lee et al. (2011) who propose seeking discriminative  $k$ -mers specific to a species not found in others, an idea implemented in CLARK (Ounit and Lonardi 2016). Another obvious approach is to represent all reference genomes at equal levels, or similarly, to sample  $k$ -mers uniformly at random. The overall goal of this work is to explore these and a set of alternative strategies in the context of taxonomic read classification and taxonomic profiling of metagenomic samples.

Subsampling  $k$ -mers from ultra-large and heterogeneous databases available these days faces many challenges. On the one hand, sampling of taxonomic groups in public libraries is very imbalanced, which can become a challenge as noted by Nasko et al. (2018). For example, RefSeq includes 35,947 *Escherichia* genomes whereas 1164 genera have only a single genome. In the data set of Zhu et al. (2019), despite their attempts to find representative species, the number of genomes sampled per taxon in each taxonomic rank varies three orders of magnitude. The true imbalance in evolutionary history (where some clades are much more diverse, abundant, and species-rich than others) in addition to biases in sampling (where some environments are sampled far more than others) are behind such imbalances. Regardless of the cause, a naive subsampling of  $k$ -mers can easily miss entire branches of the evolutionary tree. The situation is made worse by the fact that for  $k \gg 20$ , a large number of  $k$ -mers are not shared even among genomes of the same species. Examining the WoL-v1 data set, we see that 15% of pairs of genomes from the same species have pairwise distances above 5%, which results in 98% of 30-mers being unique to each genome in such pairs. Pairs of genomes from the same genus but different species shared more than 10% of their 30-mers in only 1% of cases. Thus, as the number of genomes  $N$  grows, the number of unique  $k$ -mers grows close to proportionally to  $N$ . Moreover,  $k$ -mers unique to ranks above species will be far outnumbered by those unique to a genome or a species. As a result, optimal  $k$ -mer selection is even more challenging if one tries to find a unique set of  $k$ -mers to be used at all ranks (perhaps the reason Ounit and Lonardi [2016] target a particular rank).

In this paper, we propose a set of strategies for selecting a pre-defined number of  $k$ -mers from a large pool of genomes in ways that are conducive to classifying reads. Our goal is to find subsets that do not leave out poorly sampled groups, work across taxonomic ranks, and do not reduce the ability to classify relatively novel reads. We implement and explore these strategies paired with the read classification and taxonomic profiling method CONSULT-II. Our method,  $k$ -mer RANKer (KRANK), takes a taxonomy and a set of genomes labeled with the taxonomy as input. It selects a subset of  $k$ -mers from these genomes based on its ranking strategy and a user-defined size constraint. KRANK subsamples  $k$ -mers in a bottom-up traversal of the reference taxonomy, enabling it to make choices about what  $k$ -mers to keep locally instead of globally and eliminating the need to analyze all  $k$ -mers jointly at any point in library construction. Our highly optimized and flexible C++ implementation allowed us to compare several alternative strategies empirically. Our comprehensive results on taxonomic read classification and taxonomic profiling tasks show that KRANK, unlike more naive selection strategies, can produce dramatic reductions in library sizes without substantial loss of accuracy.

## Results

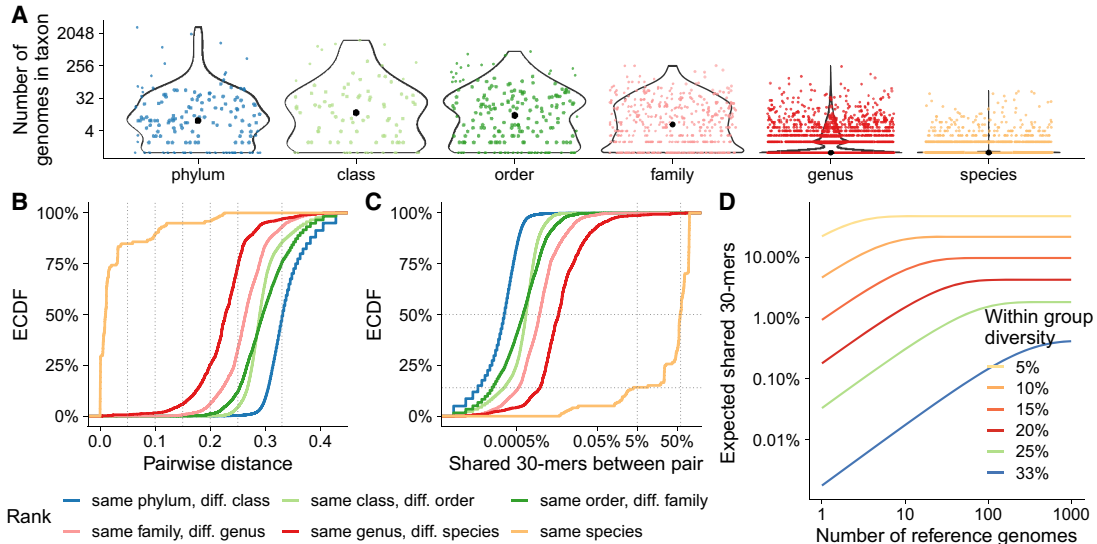
We designed KRANK (detailed in Algorithm 1) to select a pre-defined number of  $k$ -mers from a given set of genomes. KRANK designates some taxonomic rank as the leaf rank (species by default). It then traverses the taxonomy  $\mathcal{T}$  in post-order. At leaves, distinct  $k$ -mers are extracted and encoded in a compact manner (see “Algorithmic details of KRANK”). Following CONSULT-II, KRANK indexes  $k$ -mers using  $l$  locality-sensitive hashing (LSH) tables. These LSH tables ensure two  $k$ -mers with a small Hamming distance (HD; e.g.,  $\leq 3$ ) have a high chance of being indexed to the same row in at least one of the tables (see “LSH tables” for details). This is achieved by selecting  $h < k$  random but fixed positions of a  $k$ -mer as its hash key, and each  $k$ -mer is stored in the row of the table indexed by the LSH key. As we traverse the tree, on each internal node, we adaptively impose a size constraint on its tables. We combine each row of the hash tables of all the children nodes and remove  $k$ -mers based on a ranking algorithm until the size constraint is satisfied. At the root, the final  $l$  tables are restricted to  $b$  columns per row, and each will contain  $2^{hb}$   $k$ -mers. When not specified, we use  $l=2$ ,  $k=32$ ,  $h=12$ ,  $b=16$ . Given the hash tables at the root, we use the existing CONSULT-II algorithm for taxonomic classification and abundance profiling (see “CONSULT-II classification and profiling”).

The heart of the algorithm is adaptive size constraints and  $k$ -mer ranking, both of which we empirically motivate first. These explanatory sets of results use a manageable microbial data set by Zhu et al. (2019) composed of 10,575 genomes (WoL-v1 hereafter). We chose 756 query genomes in total (676 bacterial and 80 archaeal genomes), of which only 10 are present in the reference set. Measuring the novelty by the minimum distance of a query to any reference genome as estimated by Mash ( $d^*$ ), we bin queries into novelty groups (Supplemental Fig. S1A). We use 150 bp error-prone reads simulated from each genome (66,667 reads per query genome) and test methods in terms of their accuracy in classifying individual reads. See “Data sets” and “Evaluation metrics” for details. After exploratory results, we move to more formal evaluations of the method for two tasks: taxonomic classification of metagenomic reads and taxonomic profiling of metagenomic samples.

### KRANK: motivating the design

#### Adaptive size constraints

We observed wide differences among taxa, even at the same rank, in terms of sampling density (Fig. 1A) and levels of diversity (Fig. 1B). This heterogeneity creates a major challenge in selecting  $k$ -mers. If we randomly sample  $k$ -mers, taxa will contribute proportionally to their size in the final subset. Given a total budget  $M$ , in expectation, a taxon  $t$  would have  $M|\mathcal{K}_t|/|\mathcal{K}|$   $k$ -mers in the sampled subset, where  $\mathcal{K}_t$  is the set of  $k$ -mers of all genomes labeled by taxon  $t$ , and  $\mathcal{K}$  is the set of all reference  $k$ -mers. As a result, taxa with lower sampling will have little representation, whereas highly sampled groups (e.g., *Escherichia coli*) will dominate. We partially address this challenge by imposing a size constraint on hash table size for each internal node during traversal to avoid having bloated nodes that starve sister taxa. Although this balancing can in theory be done at the root, we would need to keep a map from  $k$ -mers to genomes, which is impractical. Thus, preemptively filtering some number of  $k$ -mers from some tables helps scalability by reducing the memory usage. We designed a heuristic for adaptive size constraints whereby, each node  $t$  containing portion  $r(t)$  of the total data set gets a budget of  $\sqrt{r(t)}M$   $k$ -mers. We define  $r(t)$  as either the portion of total  $k$ -mers present



**Figure 1.** Exploratory analysis of imbalanced taxon sampling and the portion of shared  $k$ -mers across taxonomic ranks in the WoL-v1 data set (Zhu et al. 2019). (A) Number of reference genomes under each taxonomic node (dots), separated by ranks. (B, C) The distribution of Mash (Ondov et al. 2016) estimated genomic distances (B) and Jaccard similarities (C) among 500,000 randomly sampled pairs of genomes that share a taxonomic rank but are different in lower ranks. The empirical cumulative distribution function (ECDF) is shown. (D) The theoretical expectation for the number of 30-mers shared between a query and at least one of  $N$  sampled genomes of a reference set for a group that has within-group diversity  $2d$  is shown as  $(1-d)^k(1-(1-d)^k)^N$ .

under  $t$  (default) or the proportion of taxonomic leaves (i.e., species) under  $t$  (see section “ADAPTIVE SIZE CONSTRAINT” for details). We empirically evaluate these two versions of this heuristic against the baseline of enforcing no adaptive constraint, all paired with random  $k$ -mer subsampling.

Imposing the adaptive size constraints can make a substantial difference for taxonomic groups with low levels of sampling (Fig. 2A). For example, for phyla with [0,50] and (50,500] reference genomes sampled, both definitions of  $r(t)$  improve the F1 accuracy over enforcing no adaptive constraint for up to 35% and 12%, respectively, with similar improvements obtained at other ranks. Thus, the adaptive size constraints achieve their stated goal of ensuring groups with lower sampling are not left out. This equitable spreading of the  $k$ -mer budget inevitably results in reduced representation for highly sampled groups. For example, for the few phyla with more than 500 reference genomes, we observe a considerable drop in the performance (e.g., F1 score declines by 14%–20%). Thus, it may appear that better sampling of low-representation groups has to come at the expense of the high-representation groups. However, this trade-off can be avoided.

Because we have two different tables ( $l=2$ ) in the library, we can have a mixed strategy: We enforce the adaptive size constraint for one table but not the other. This mixed strategy is almost as good as the adaptive version for low-representation groups, and almost as good as the no-adaptive-constraint strategy for highly sampled groups (Fig. 2A). Overall, this mixed strategy outperforms having no adaptive constraints, both in terms of precision and recall (Fig. 2B). The improvements are most visible for moderately novel queries ( $d^* \in (0.025, 0.1]$ ). Across queries, the mixed strategy improves the recall by 6%–10% and precision by 2%–5% above the genus rank. Thus, KRANK adopts the mixed strategy as the default.

### Ranking and removing $k$ -mers

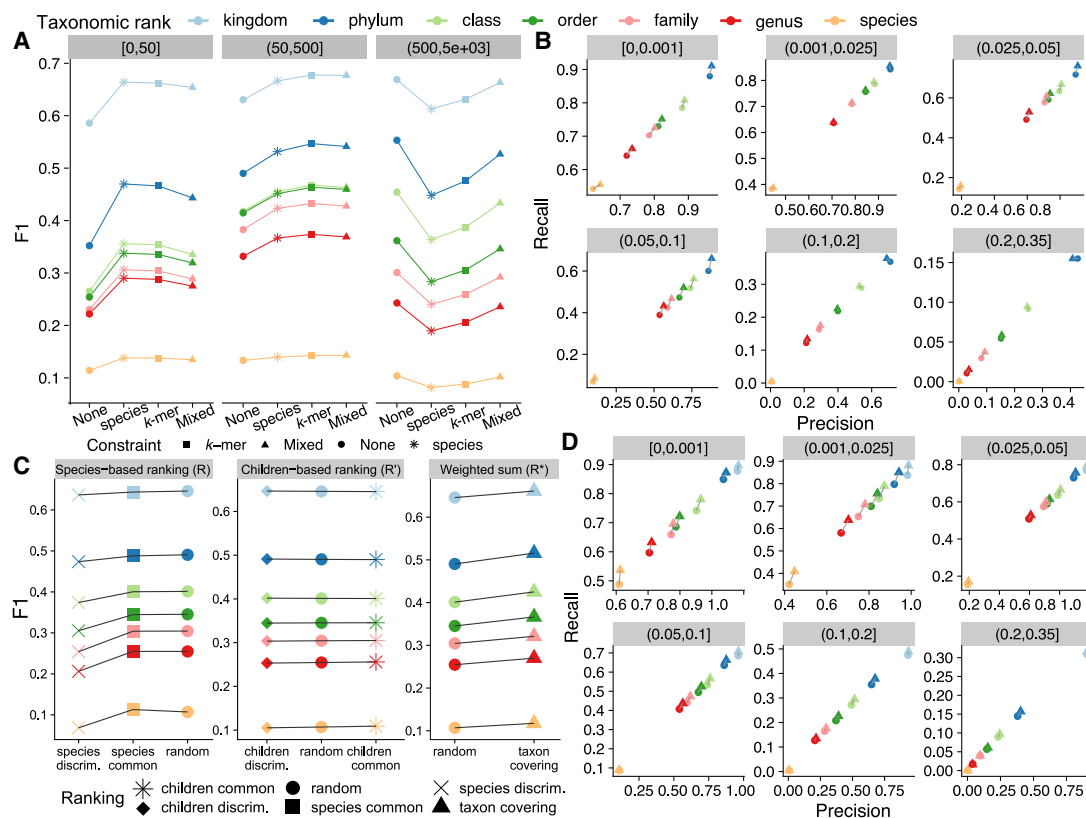
Faced with the question of which  $k$ -mers to select given the limited budget, three options immediately present themselves: selecting

randomly, selecting *discriminative*  $k$ -mers unique to some taxa, or conversely, selecting common  $k$ -mers present in many taxa. We compared these options and others that we propose paired with the adaptive size-constraint approach described above using constraints for both tables (i.e., not mixed). In implementing these strategies, we used two quantities—the number of leaves (i.e., species) under taxon  $t$  that include a  $k$ -mer  $x$  denoted by  $R(x,t)$ , and the number of children of  $t$  that include the  $k$ -mer  $x$  at least once denoted by  $R'(x,t)$ , as shown in Supplemental Figure S2.

### The case against discriminative $k$ -mers

Keeping discriminative  $k$ -mers found only in specific target species, an approach followed by Ounit et al. (2015), would correspond to preferentially filtering out  $k$ -mers that appear in more than one species (or some other rank), which corresponds to  $R(x,t) > 1$ . Going one step further, we emulate the discriminative  $k$ -mer strategy by ranking  $k$ -mers inversely to  $R(x,t)$  or  $R'(x,t)$ , thus removing common  $k$ -mers. Empirically, using either version is worse than simply using random selection (Fig. 2C). The reduction is negligible with children-based  $R'(x,t)$  but is quite dramatic when used with species counts  $R(x,t)$ ; for example, the declines in F1 are 36% and 19% for species and genus ranks, respectively. Thus, we caution against using discriminative  $k$ -mers in a memory-bound setting.

The reduced accuracy of selecting unique  $k$ -mers can be explained by noting that very few  $k$ -mers are shared even among different species of the same genus. For example, in the WoL-v1 data set (Zhu et al. 2019), a random sample of genomes that are from different species but the same genus share only 0.05% of their 30-mers (Fig. 1C) and have Mash distances that typically exceed 20% (Fig. 1B). With simplifying assumptions, we can approximate the probability that a  $k$ -mer would be shared between a query genome and at least one of  $N$  genomes sampled from a taxonomic group, assuming each pair of genomes in this group has distance  $d$  to parent, as  $(1-d)^k(1-(1-d)^k)^N$ . Plotting this equation



**Figure 2.** Evaluation of read classification performance of KRANK using different heuristics for  $k$ -mer selection. (A,B) Comparison of adaptive size-constraint approaches with random  $k$ -mer selection. None: enforcing no constraint, species: number of species constraint,  $k$ -mer: total  $k$ -mer count constraint, Mixed: no constraint for one table and  $k$ -mer count constraint for the other. (A) F1 accuracy, dividing queries (panels) into three bins based on the number of reference genomes in the phylum of the query. (B) Precision versus recall of None versus Mixed (default), dividing queries into bins of novelty ( $d^*$ ). (C,D) Evaluation of simple ranking strategies (discriminative, common, random) implemented using species counts  $R$  or children counts,  $R'$  and our new weighted-sum,  $R^*$  heuristic [Equation (11)]. See Supplemental Figure S2 for an illustration of functions. (C) F1 across all queries. (D) Precision versus recall of random and weighted-sum heuristic, dividing queries by bins of novelty ( $d^*$ ). In (A) and (C),  $x$ -axes are ordered by the average F1 score of all queries. All results use:  $w = 35$ ,  $k = 32$ ,  $l = 2$ ,  $h = 12$ , and  $b = 16$ .  $d^*$  is the distance to the closest reference estimated by Mash. Default KRANK uses mixed adaptive size and  $R^*$  ranking.

shows that using discriminative  $k$ -mers leads to very few matches between a query and the reference (Fig. 1D). For instance, in a group with  $2d = 20\%$  diversity (which is less than most genera; see Fig. 1B), only 0.7% of query 30-mers are expected to be found in at least one among  $N = 5$  reference genomes and only 4.2% when  $N \rightarrow \infty$  (infinitely many genomes sampled from that genus). Because most  $k$ -mers are unique and subsampling is inevitable in a memory-bound setting, removing common  $k$ -mers makes it difficult to find  $k$ -mer matches between reads and reference species that have 20% or more distance to them (e.g., when the closest available match is at the genus level). This only gets worse at higher ranks and is not much better within the same species, where the pairwise Mash distance is 5% or more in 15% of cases.

#### Selecting common $k$ -mers alone does not help

We can emulate common  $k$ -mer selection by ranking  $k$ -mers proportionally to  $R(x,t)$  or  $R'(x,t)$ . Given the argument against discriminative  $k$ -mers, common  $k$ -mers seem appealing because they can represent each taxon using conserved sequences present across different genomes. Empirically, this strategy used with either  $R(x,t)$  or  $R'(x,t)$  is better than using discriminative  $k$ -mers but is no better than random sampling (Fig. 2C). Just like randomly selecting  $k$ -mers, this strategy focuses the budget on larger and densely sam-

pled taxa. For instance, the WoL-v1 reference set includes 2975 *Pseudomonadota* genomes and only a single genome from many other phyla. Although some unevenness is undoubtedly due to genuine differences among the diversity of phyla, some must be due to uneven sampling.

#### Covering all children improves accuracy

Instead of maximizing total coverage, we can attempt to cover all species. We designed a scalable heuristic ranking mechanism to fulfill this goal. We rank  $k$ -mers by a weighted-sum of  $R(x,.)$  values of  $t$ 's children (see Method section "FILTERBYRANK" and  $R^*(x,t)$  defined in Equation (11) and illustrated in Supplemental Fig. S2). The weights are used to de-emphasize children that are highly sampled among surviving  $k$ -mers. We simply set the weights to be inversely proportional to the coverage of each child taxon among  $k$ -mers of a particular row of the LSH table to ensure that children covered with fewer surviving  $k$ -mers get covered by the parent. Using this approach improves  $k$ -mer selection consistently across all taxonomic ranks compared to the random strategy and common  $k$ -mer strategy (Fig. 2C). Further dividing queries by their novelty, we observe that our weighted-sum strategy improves both precision and recall, particularly for less novel queries; recall

increases by as much as 15% at the species level in the  $d^* \in (0.001, 0.25]$  bin (Fig. 2D).

### Comparison to other methods for taxonomic read classification

The default version of KRANK uses both adaptive size constraints (mixed strategy) and *k*-mer ranking and comes in two default modes: high-sensitivity (hs) using 51.2 GB and lightweight (lw) using 12.8 GB. We next compare these two versions of KRANK against CONSULT-II, CLARK, and Kraken 2 on the same WoL-v1 data set used so far in terms of accuracy in classifying individual reads (Fig. 3).

#### Accuracy in default settings

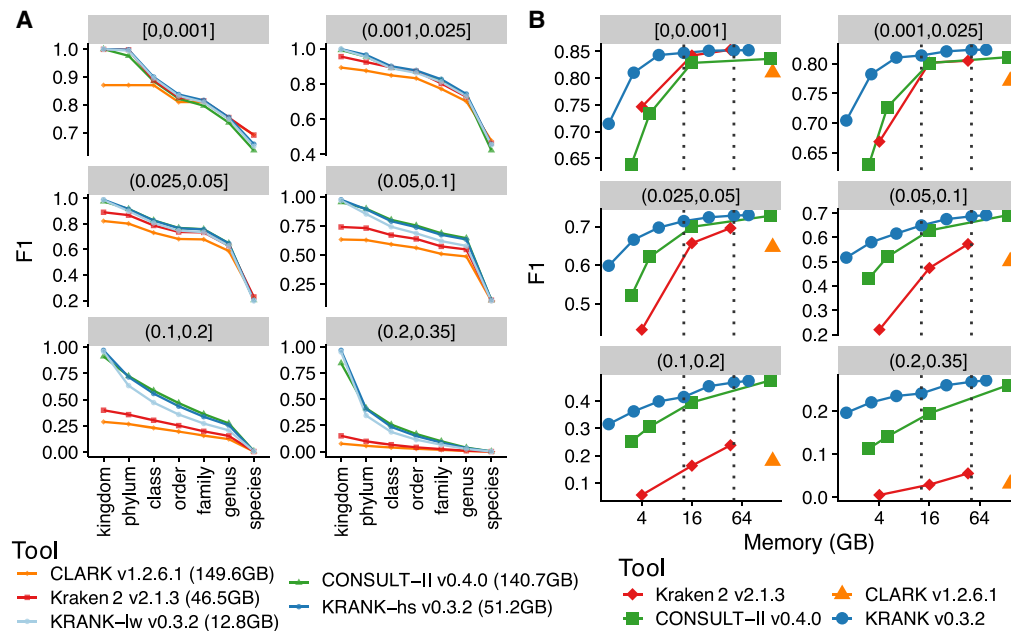
Across different ranks and query novelty levels, KRANK-hs matched or improved on CONSULT-II in terms of F1 accuracy despite using far fewer *k*-mers and about one-third of the memory (Fig. 3A). KRANK-hs and CONSULT-II often had similar precision (with a slight advantage for CONSULT-II in some cases), but KRANK-hs had better recall in many settings (Supplemental Fig. S3). Similar to CONSULT-II, KRANK achieved substantially higher recall compared to Kraken 2 and CLARK. Improvement in recall usually comes with little or no sacrifice in precision, which can be controlled with the total vote threshold parameter of CONSULT-II (Şapci et al. 2024). KRANK-lw is more accurate than CONSULT-II at the kingdom level and only slightly less accurate elsewhere despite using 10 times less memory. Accuracy also depended on query novelty ( $d^*$ ). All methods performed similarly both in terms of precision and recall for queries that were not novel ( $d^* \leq 0.025$ ), except CLARK which had slightly lower accuracy (Supplemental Fig. S3). For these less novel queries, both variants of KRANK were indistinguishable from CONSULT-II. For more

novel queries ( $d^* > 0.025$ ), accuracy dropped for all methods and substantial differences between them emerged. For example, for the  $d^* \in (0.05, 0.1]$  bin, KRANK-lw outperformed Kraken 2 and CLARK above the species rank, achieving 16% and 35% higher F1 scores, respectively, at the phylum level. Unlike KRANK-hs, KRANK-lw had slightly lower accuracy than CONSULT-II at some ranks (e.g., 8%–9% at the order rank) for  $d^* \in (0.025, 0.1]$  and moderately lower accuracy (e.g., 12% at the phylum rank) for  $d^* > 0.1$ . Nevertheless, KRANK-lw outperforms Kraken 2 (e.g., 77% and 250% higher F1 at the phylum rank for  $d^* \in (0.1, 0.2]$  and  $(0.2, 0.35]$ , respectively).

Both KRANK-hs and KRANK-lw are slightly better than CONSULT-II in terms of kingdom-level classification for novel genomes. This improvement can be due to the better representation of archaea with careful *k*-mer selection. Note that our query set is highly imbalanced (e.g., 90% bacteria and 33% *Pseudomonadota*). Thus, simply assigning all reads to the largest group (in terms of the number of query genomes) at each rank would still perform well (Supplemental Fig. S4). As expected, this null model performs well at the kingdom level (0.94 F1 score). However, KRANK (both lw and hs versions) is the only method that improves on the null model (0.96 and 0.97 average F1 scores, respectively) and assigns 83% of archaeal reads to archaea. At lower ranks, all tools, except for CLARK at the phylum rank, do substantially better than the null model.

#### Adjusting memory usage

We next explored three to seven levels of consumed memory for KRANK, CONSULT-II, and Kraken 2 to compare their accuracy given similar levels of memory, ranging from 1.6 GB to 140 GB (Fig. 3B). KRANK clearly outperformed both alternatives in terms of



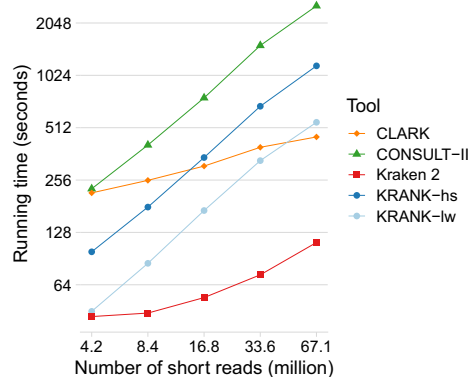
**Figure 3.** Detailed evaluation of read classification performances of tools across different ranks and novelty levels. (A) Comparison of two default modes of KRANK (lw for lightweight and hs for high-sensitivity) with other methods across ranks and novelty bins on the WoL-v1 data set with 66,667 simulated 150 bp reads from each of 756 query genomes. Each panel is a novelty bin, corresponding to the minimum distance to any reference genome ( $d^*$ ), measured by Mash (Ondov et al. 2016); see Supplemental Figure S1. F1 scores are averages of all queries in each  $d^*$  bin. See Supplemental Figure S3 for precision and recall. (B) F1 accuracy for KRANK, CONSULT-II, and Kraken 2 as we change the memory used in GB. CLARK is shown as a single data point. We mark memory levels corresponding to the lw and hs modes using dashed lines.

average F1 score when memory is controlled. At lower memory levels, CONSULT-II had much lower F1 scores than KRANK, and Kraken 2 had lower F1 than CONSULT-II for more novel queries. For instance, given  $\leq 4$  GB, KRANK was 26% more accurate than CONSULT-II and 8% more accurate than Kraken 2 for the least novel bin. For the moderately novel queries ( $d^* \in (0.05, 0.1]$ ), KRANK with 3.2 GB outperformed Kraken 2 with 46.5 GB. KRANK with 1.6 GB outperformed CLARK (149.6 GB) and Kraken 2 (46.5 GB) for novel queries ( $d^* > 0.1$ ), giving 160% and 89% higher F1 scores, respectively. Note that CONSULT-II and KRANK store the same number of  $k$ -mers (268 M in a single table) when memory usages are 5 GB and 3.2 GB, respectively. KRANK remained more accurate than CONSULT-II when we compared the method for varying numbers of  $k$ -mers kept in the library (Supplemental Figs. S5 and S6). Because the only difference between KRANK and CONSULT-II is the selected  $k$ -mers, these results demonstrate the effectiveness of the selection algorithm.

### Resource usages

Despite performing more complex computations, KRANK built libraries using less computational resources due to its efficient implementation, compared to CONSULT-II (Table 1). Although the total time needed was also less, the main gain was KRANK's distributed memory implementation. KRANK splits the LSH table into batches, enabling it to divide the workload into independent jobs. For instance, on the WoL-v1 data set, KRANK-hs was built using 512 batches per table, each of which took 8 min on average with four threads. KRANK-lw used 256 batches and built each batch in 5 min on average. Together with the initial preprocessing of reference genomes, processing 32 batches in parallel (e.g., on different cluster nodes) needed 4.5 h for high-sensitivity mode to construct a library with two tables. Using this configuration, the peak memory usage for each batch was  $< 32$  GB. Although the running times are not directly comparable because of the difference in parallelism, KRANK needed significantly shorter times than CONSULT-II and CLARK to build its libraries.

At the query time, compared to CONSULT-II, KRANK does not introduce any extra computation but its smaller libraries resulted in reduced running times. On this data set, KRANK-hs was more than twice as fast as CONSULT-II (Table 1). These improvements are due to better cache performance and shorter library loading times. Note, however, that Kraken 2 was  $\sim 13\times$  faster than KRANK-hs and scaled better as the number of queried short reads increased (Fig. 4). This difference may be due to the fact that KRANK needs to compute HD for up to  $2b=32$  reference



**Figure 4.** Query time for varying numbers of short reads sampled across 756 genomes from different novelty levels. Running time measurements were performed using 96 threads on a machine with 2.2 GHz AMD EPYC 7742 processors. Both axes (x and y) are in log-scale.

$k$ -mers for each query  $k$ -mer, whereas Kraken 2 only computes presence/absence.

### Abundance profiling: CAMI challenges

#### CAMI-I high-complexity data set

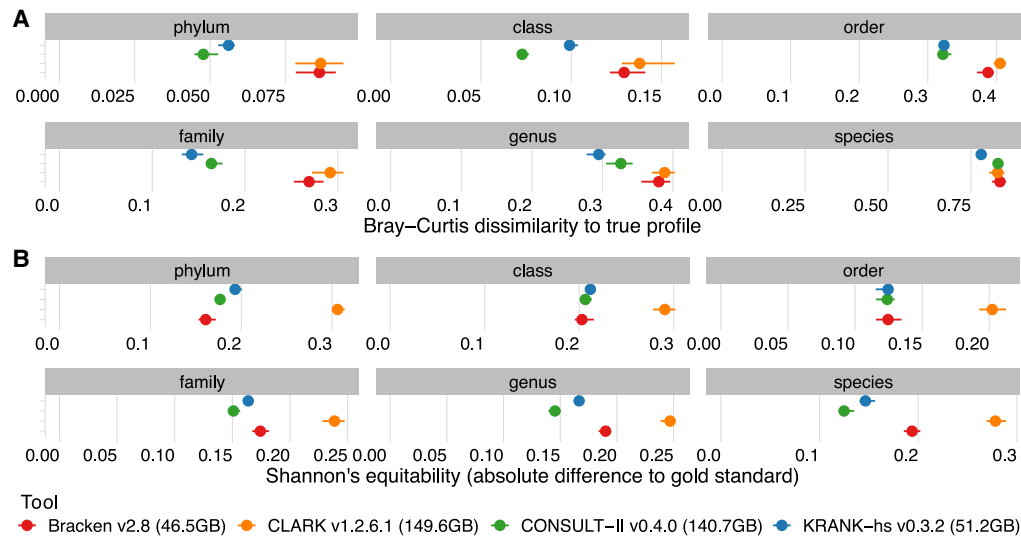
We next compared KRANK-hs against three alternatives in terms of taxonomic profiling accuracy on the Critical Assessment of Metagenome Interpretation (CAMI-I) (Sczyrba et al. 2017) benchmarking challenge, focusing on the high-complexity subset (can be found at <http://gigadb.org/dataset/100344>). We compared against CONSULT-II, CLARK, and Bracken (a Bayesian extension of Kraken 2) using standard metrics promoted by CAMI-I (see “Data sets” and “Evaluation metrics”). All methods were run with the custom libraries built using the same WoL-v1 data set used in the previous analysis.

Abundance profiles estimated by KRANK-hs and CONSULT-II are the most similar to the true profile in terms of Bray–Curtis dissimilarity (Fig. 5A). At the species level, all methods have high errors and are comparable, with a slight advantage for KRANK-hs. As we move up the ranks, the error quickly drops for all methods, but KRANK-hs and CONSULT-II outperform others. Overall, KRANK-hs performs similarly to CONSULT-II despite the much-reduced memory. KRANK performs considerably better than both Bracken and CLARK. All methods tend to underestimate Shannon's equitability, which is a measure of the variety and distribution of taxa

**Table 1.** Computational resources needed for a database built on 10,470 reference genomes and 756 queries (50 M short reads)

	Running time		Peak memory	
	Query	Library	Query	Library
Kraken 2	88 sec, 96 threads	3.5 h, 128 threads	47 GB	50 GB
CLARK	676 sec, 96 threads	13.5 h, 128 threads	150 GB	370 GB
CONSULT-II	2191 sec, 96 threads	19 h, 128 threads	141 GB	157 GB
KRANK-hs	1167 sec, 96 threads	4.5 h, $32 \times 4$ threads	51 GB	32 GB per batch
KRANK-lw	787 sec, 96 threads	2.5 h, $32 \times 4$ threads	13 GB	8 GB per batch

Measurements for both queries and library building were performed on a machine with 2.2 GHz AMD EPYC 7742 processors. Reported library building times are for 256 and 512 batches, respectively, for KRANK-lw and KRANK-hs, each with four threads, distributed across 32 cluster nodes and run in parallel.



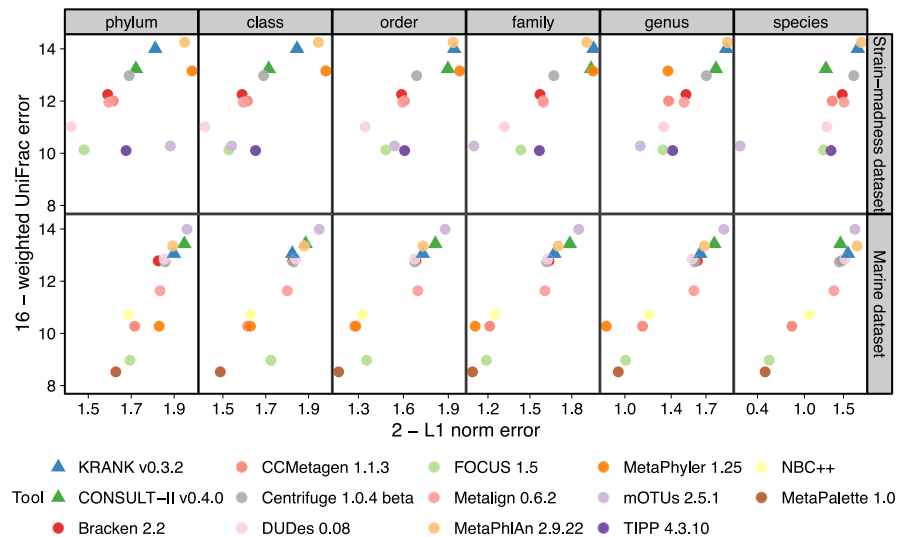
**Figure 5.** Taxonomic profiling on CAMI challenge. (A) The Bray–Curtis metric measures the dissimilarity between the estimated and true profiles. (B) The difference between Shannon’s equitability of the estimated profile and of the gold standard, measuring how well each tool reflects the alpha diversity (i.e., the evenness of taxon abundances). Reported metrics are based on the profile estimates of the high-sensitivity (hs) setting of KRANK. We show the mean and minimum/maximum across five samples.

present in a sample. At the phylum rank, Bracken has closer values to the gold standard, whereas CONSULT-II and KRANK become better at the family rank or lower (Fig. 5B). Overall, KRANK-hs is slightly less accurate than CONSULT-II in terms of this metric across ranks. Note that correcting for genome size improves the accuracy of both CONSULT-II and KRANK, but KRANK benefits slightly more in terms of the Bray–Curtis metric (Supplemental Fig. S7A). Because the largest three phyla in our reference set constitute 80%–85% of each sample, this data set presents a case where the KRANK strategy of covering low-sampled groups is not needed and sampling of highly covered phyla by CONSULT-II is sufficient.

#### CAMI-II marine and strain-madness data sets

To compare KRANK to a larger set of methods, we built a library based on 72,766 genomes available on RefSeq as of January 8, 2019. Considering the high number of reference genomes, we omitted the lightweight memory configuration and focused on the high-sensitivity setting, which already uses less than half of the memory used by CONSULT-II. Because this same data set had been provided to all methods tested in the CAMI-II challenge (Meyer et al. 2022), we were able to compare KRANK with 12 alternative methods from that original paper, in addition to CONSULT-II, which was also applied to this data set (see “Data sets”). These methods include both marker-based and *k*-mer-based methods, which use reads from the entire genome as opposed to relying on predefined markers.

Taxonomic profiles of KRANK-hs were among the most accurate on the CAMI-II challenge (Fig. 6). On the strain-madness data set, in which only two phyla constitute 90% of the relative abundance, our method was a close second to the best method, which was the marker-based method MetaPhlAn (Fig. 6; Supplemental Fig. S8A). Averaged across all samples, the UniFrac score of KRANK-hs was only 2% less than MetaPhlAn and 6% above the third-best method, CONSULT-II. In species, class, and phylum ranks, MetaPhlAn had slightly better L1 than KRANK-hs, whereas



**Figure 6.** Comparing KRANK (high-sensitivity setting using 51.2 GB) with other participants in CAMI-II benchmarking challenge (Meyer et al. 2022). As in that original paper, we show the upper bound of L1 norm (2) minus actual L1 norm versus the upper bound of weighted UniFrac error (16) minus actual weighted UniFrac error. Each data point stands for the average of 100 strain-madness samples or 10 marine samples. Metrics were computed using OPAL with default settings and the `-n` option. The UniFrac error is the total amount of predicted abundances that needs to be moved along the edges of the taxonomic tree to make them overlap with the true abundance profile. The L1 error simply measures the accuracy of reconstructing the relative abundance profile at a fixed rank.



in the family rank, KRANK-hs was slightly better (Fig. 6). Note that KRANK and MetaPhlAn are based on very different approaches, and the reason behind MetaPhlAn performing better is not immediately clear. It may be that using markers performs better than using entire genomes in some cases as markers are less subject to horizontal gene transfer. Because MetaPhlAn lacks the capability of classifying individual reads, we were not able to include it in the read classification experiments. Despite using the same profiling algorithm as CONSULT-II and one-third of its memory, KRANK-hs also had better L1 metrics than CONSULT-II across all ranks on this data set. Other *k*-mer-based methods such as Bracken had far lower accuracy in terms of both metrics.

On the marine data set, which includes both archaeal and bacterial genomes, and consists of 65%–70% *Proteobacteria*, the pattern slightly changed (Fig. 6). Here, CONSULT-II was slightly more accurate than KRANK-hs, with 3% better UniFrac and 3% better L1, averaged over all ranks and samples. Only at the species level, KRANK-hs was more accurate than CONSULT-II (6% in terms of L1). Overall, KRANK-hs was ranked fourth, following mOTU, CONSULT-II, and MetaPhlAn (Supplemental Fig. S8). Note, however, that the best method, mOTU, performed poorly on the strain-madness data set. KRANK-hs was the second-best *k*-mer-based method after CONSULT-II but used only one-third of its memory. The third leading *k*-mer-based method, Bracken, was less accurate than CONSULT-II and KRANK-hs in both metrics, but with a smaller margin than the strain-madness data set.

To summarize, KRANK-hs was the best *k*-mer-based method and the second-best method overall (Supplemental Fig. S8B). Averaged over both data sets, the marker-based method MetaPhlAn had the best accuracy, followed closely by KRANK-hs, and then CONSULT-II. Note, however, that we focused on the two main CAMI-II metrics that accounted for abundance values. Judged by the presence/absence of taxa (i.e., ignoring abundance), KRANK-hs had better purity than CONSULT-II and Bracken and about the same level of completeness (Supplemental Fig. S9). However, these abundance-agnostic metrics were not the strength of *k*-mer-based methods; across both data sets, KRANK-hs, CONSULT-II, and Bracken were among the best in terms of completeness but not purity.

## Discussion

We introduced KRANK for selecting a memory-bound and representative subset of *k*-mers from a large genome collection for the purpose of taxonomic classification and abundance profiling. We used a tree-based *k*-mer filtering algorithm that decides which *k*-mers to remove as it moves up a taxonomic tree. This bottom-up strategy makes the library construction more memory efficient as *k*-mers are analyzed in smaller chunks at each node. In addition, we use an adaptive size constraint mechanism to leave some filtering decisions to higher taxonomic levels. These decisions are made based on the prevalence of *k*-mers among children of or leaves under each node.

Ultimately, both our adaptive size constraints and our ranking strategies are heuristics, with no theoretical guarantees, but designed based on careful experiments involving novel queries and imbalanced sampling. Due to the empirical nature of the method, the exploration of alternatives will be an ongoing topic of research, one that is helped by KRANK's software design. In particular, our approach still leaves room for ties during ranking, and smarter ways to break such ties may help accuracy. Alternatively, more theoretically justified approaches should be explored.

One obstacle to building a more theoretical approach is to define an objective function for *k*-mer selection. Subtleties of imbalanced and biased sampling in reference databases make the definition of the desired optimization problem tricky. Should we attempt to combat the sampling bias by selecting *k*-mers more proportionally from under-sampled groups? If so, how can we know and model the bias? Should we attempt to consider likely profiles of queries? If so, under what model? These questions have no easy answers. Nevertheless, defining mathematically rigorous but nuanced objectives for *k*-mer selection could be a fruitful avenue of future research.

Our approach focused on *k*-mer similarity along the evolutionary dimension but completely ignored *k*-mer overlap. Selecting *k*-mers based on overlap is well covered by minimizers, which we use as a preprocessing step. However, our ranking algorithm is unaware of overlap and minimizers are unaware of the taxonomic relationships. Integrating the evolutionary-aware strategies of KRANK with location-based minimizers poses an interesting algorithmic challenge and could further improve the method. However, designing such a combination is not trivial given our hashing strategy and needs further algorithmic developments.

Despite all the improvements compared to CONSULT-II, the library construction of KRANK is expensive in terms of CPU hours spent. Nevertheless, the creation of the reference library is a one-time operation, amortized over many subsequent uses. Furthermore, computation is getting cheaper rapidly, and the total resources we used (e.g., <550 CPU hours for the WoL-v1 library) are far from prohibitive. Crucially, the batching strategy enables us to build the library over a cluster (i.e., distributed memory), with each job using as little as 8 GB and many running in parallel in the order of minutes (Table 1). Moreover, note that the library building time is a function of both the number of genomes and the size of the taxonomy, which grows less rapidly. For example, whereas the CAMI-II library had 7× more genomes than WoL-v1, we needed only 3×, because the taxonomy had 1.8× as many nodes, and 2× as many species as WoL-v1. Our analyses of 72,766 genomes included in the CAMI-II case demonstrate the scalability of the method to scales commensurate with modern-day libraries.

As data sets become even larger, KRANK parameters may need to adjust. Many ultra-large reference databases are available now, including GTDB (Parks et al. 2018) with 402,709 genomes, and WoL-v2 (Balaban et al. 2024), with a phylogeny built from 199,330 genomes. Moreover, RefSeq currently includes millions of genomes. The fast increase is primarily due to the sequencing of many more genomes from key species. As the number of genomes from each species increases, our choice in this paper to treat the species rank as the leaves of the taxonomy becomes less scalable. However, KRANK can easily use a more refined group (e.g., subspecies) as the leaf, which will trade-off slightly longer running times for better memory efficiency. It could also enable subspecies-level classification, a topic that should be explored in the future. Another concern in handling ever-growing data sets is the need for de novo construction of the database as new genomes become available. Future research should explore ways to add new genomes into existing libraries, without repeating all the steps.

## Methods

As a *k*-mer selection method, KRANK can, in principle, be paired with any *k*-mer-based classification or profiling methods. However, in this paper, we use KRANK solely with CONSULT-II (Şapcı et al. 2024); *k*-mer matches found in libraries built by

KRANK are given as input to CONSULT-II v0.5.0 for classification and profiling. Moreover, our  $k$ -mer selection algorithm adopts a similar LSH data structure to CONSULT-II. Thus, we start by describing the CONSULT-II hash tables and the CONSULT-II classification and profiling methods before presenting KRANK.

### LSH tables

CONSULT(-II) was designed to answer the following question: Are there any  $k$ -mers in a given reference set with HDs less than some threshold  $p$  to a given query  $k$ -mer? The core idea is to find low HD matches efficiently by partitioning reference  $k$ -mers to constant-sized subsets (i.e., rows of the hash table with size  $b$ ) using the bit-sampling LSH method (Har-Peled et al. 2012). In this method, we select  $h$  random but fixed positions of a  $k$ -mer as its hash key. Each LSH table is a simple  $2^{2h} \times b$  array, and  $l$  such hash tables constitute the data structure (default:  $l=2$ ,  $h=14$ , and  $b=16$  for KRANK).

Each query is compared against all  $l \times b$   $k$ -mers with the same hash key. We return the closest  $k$ -mer match with distance  $\leq p$  if any exists. Because we explicitly compute the HD of the query to these reference  $k$ -mers, there are no false positive matches but false negatives are possible. Using the bit-sampling scheme (Har-Peled et al. 2012), one can guarantee that two  $k$ -mers at HD =  $d$  have the same hash with probability  $(1-d/k)^h$  assuming independence of positions. Thus, given two  $k$ -mers, the probability that *at least one* of  $l$  LSH keys is the same for both  $k$ -mers is given by

$$\rho(d) = 1 - \left(1 - \left(1 - \frac{d}{k}\right)^h\right)^l. \quad (1)$$

For small enough  $p$  (e.g.,  $p=4$ ), the probability  $\rho(d)$  drops quickly when  $d \gg p$  and it is sufficiently high if  $d \leq p$ . Furthermore, because classification is done at the read level, there will be  $L-k+1$  chances for a  $k$ -mer match for a read of length  $L$ . For such a read, the expected number of matching  $k$ -mers is  $(L-k+1)\rho(d)$ , which can be large enough for several realistic choices of  $l$  and  $h$  (Supplemental Fig. S10). For instance, with  $k=30$ ,  $h=14$ , and  $l=2$  (KRANK-hs settings), for a 150 bp read at 15% distance from the closest reference genome, we would still expect 23.5  $k$ -mer matches *if all of those reference  $k$ -mers are in the database*. If only 20% of the  $k$ -mers from the reference genome are in the library, we still expect more than 5  $k$ -mer matches, which can be sufficient for classification. In practice, the number of matches will depend on the budget  $M$ , and will necessarily decrease as we reduce  $M$ .

CONSULT-II extends LSH tables further by adding a taxonomic label for each  $k$ -mer in the table. This label is the lowest common ancestor (LCA) of a *subset* of species that have the corresponding  $k$ -mer (hence the name *soft-LCA*). This subset is determined with a probabilistic procedure that performs a Bernoulli trial for each genome and includes its species in the subset if the trial is successful. The success probability is different for each  $k$ -mer  $x_i$  and is a function of the number of genomes in which the  $k$ -mer appears ( $N_i$ ) set to  $p_u(N_i) = w \log_2^{-1}(2^{(N_i-1)/s} + 2)$ , where  $w$  and  $s$  are hyperparameters.

### CONSULT-II classification and profiling

Both read classification and profiling algorithms of CONSULT-II are based on voting. Given a query read  $\mathcal{R}$ , CONSULT-II lets each  $k$ -mer  $x \in \mathcal{R}$  vote for the soft-LCA of the reference  $k$ -mer it match-

es. The value of the vote decreases as the HD between the query  $k$ -mer  $x$  and its best match increases. Let the best match be at distance  $d$  to a  $k$ -mer with soft-LCA at taxon  $t$ ; then,  $x$  votes  $v_x(t) = (1-d/k)^k$  units for  $t$ . To take hierarchical relations in the taxonomic tree into account, CONSULT-II aggregates votes by recursively summing them in a bottom-up manner. Let  $\text{child}(t)$  be the set of children of the taxon  $t$ . Then, the total vote of taxon  $t$  for read  $\mathcal{R}$  is

$$\bar{v}_{\mathcal{R}}(t) = \sum_{x \in \mathcal{R}} v_x(t) + \sum_{t' \in \text{child}(t)} \bar{v}_{\mathcal{R}}(t'). \quad (2)$$

For read classification, CONSULT-II simply determines a majority vote threshold  $\tau$ , corresponding to half of the total vote value at the root of the tree. It then assigns each read to the group at the lowest rank whose total vote value strictly exceeds  $\tau$ . Note that this choice is guaranteed to be unique.

To derive an abundance profile, CONSULT-II first normalizes the total vote values of a read  $\mathcal{R}_i$  per each taxonomic rank using

$$v_{\mathcal{R}_i}^*(t) = \frac{\bar{v}_{\mathcal{R}_i}(t)}{\sum_{t' \in \mathcal{T}_r} \bar{v}_{\mathcal{R}_i}(t')}, \quad (3)$$

where  $\mathcal{T}_r$  is the set of all taxa at rank  $r$  (e.g., all phyla). Next, it gathers normalized total vote values of all  $n$  reads  $\mathcal{R}_1, \dots, \mathcal{R}_n$  in a sample, and normalizes again to obtain the final profile. Let  $p^r = [p_t^r]_{t \in \mathcal{T}_r}$  denote the relative read abundance profile at rank  $r$ , summing up to 1. CONSULT-II (since v0.3.0) sets the relative abundance of taxon  $t$  to:

$$p_t^r = \frac{\sum_i v_{\mathcal{R}_i}^*(t)}{\sum_{t' \in \mathcal{T}_r} \sum_i v_{\mathcal{R}_i}^*(t')}. \quad (4)$$

Note that  $p_t^r$  estimates the ratio of *reads* belonging to the taxon  $t$ , but we are often interested in the relative abundances of *cells* (denoted by  $\hat{p}^r = [\hat{p}_t^r]_{t \in \mathcal{T}_r}$ ). Thus, CONSULT-II (since v0.4.0) incorporates genome sizes by simply correcting  $p_t^r$  values as follows:

$$\hat{p}_t^r = \frac{p_t^r l_t^{-1}}{\sum_{t' \in \mathcal{T}_r} p_{t'}^r l_{t'}^{-1}} \quad (5)$$

where  $l_t$  is the average genome length of all references in taxon  $t$ . See Supplemental Figure S7A for a comparison between v0.3.0 corresponding to Equation (4) and v0.4.0 implementing Equation (5) on the CAMI-I data set. As an option, CONSULT-II can replace the denominator of Equation (4) with the total vote at the root of the taxonomic tree to quantify the abundance of unclassified taxa. This can be viewed as propagating vote values down from a parent to an artificial "other" lineage which continues until the species rank.

As opposed to classification, CONSULT-II's profiling (v0.4.0) algorithm does not require a majority vote, and hence, a read can contribute to abundance values of multiple taxa at a fixed rank, even when there exist only a few high HD matches. Furthermore, normalized values can accumulate over billions of reads in a sample, and can result in erroneous profiles (potentially explaining relatively high errors of CONSULT-II at species rank in CAMI-II strain-madness data set). In v0.5.0 (first introduced in the present manuscript), CONSULT-II introduces a new profiling algorithm, which uses a very similar principle to its classification algorithm. Keeping the genome size correction (Equation (5)) the same, v0.5.0 modifies Equation (4) as follows:

$$p_t^r = \frac{1}{n} \sum_i \mathbb{1}\{v_{\mathcal{R}_i}^*(t) > \max(\tau, \tau_0)\} \quad (6)$$

where  $\tau$  is the total vote value of the taxonomic tree and  $\tau_0$  is a hyperparameter hard threshold (default: 0.03). In words, instead of each  $k$ -mer of a read voting for a different taxon, each read votes collectively; the read votes for the same taxon where CONSULT-II would classify it if its total vote is high enough (i.e.,  $>\tau_0$ ). Reads with total votes below  $\tau_0$  and few or ambiguous matches are considered unclassified. Thus, the resulting abundance profile does not have to sum to 1. Using v0.5.0 instead of v0.4.0 with KRANK's library resulted in improvements in terms of Bray-Curtis dissimilarity especially at lower ranks (Supplemental Fig. S7). Thus, we report profiling results using v0.5.0 for KRANK in both CAMI challenges.

### Algorithmic details of KRANK

We use  $\mathcal{T}$  to denote the set of all nodes of the taxonomy when clear by context. We let  $\mathcal{S} \subseteq \mathcal{T}$  be the set of species, and for  $t \in \mathcal{T}$ , we let  $\mathcal{S}_t \subseteq \mathcal{S}$  be the set of species under node  $t$ . We can choose the leaves of the tree  $\mathcal{T}$  to correspond to any desired rank, but by default, we use species ( $\mathcal{S}$ ) as the lowest rank; this can be easily changed to any other rank. Let  $\mathcal{K}_t$  denote the set of  $k$ -mers of all genomes labeled by a taxon  $t \in \mathcal{T}$ . Notice that  $\mathcal{K}_t = \cup_{t' \in \mathcal{S}_t} \mathcal{K}_{t'}$ .

Recall that during the tree traversal of Algorithm 1, for each node  $t \in \mathcal{T}$ , three key operations are performed to build a new hash table, denoted by  $\mathbf{H}_t$ :

1. UNIONTABLES: takes the union of hash tables of its children  $\mathbf{H}_{t'}$ ,  $t' \in \text{child}(t)$ ,
2. ADAPTIVESIZECONSTRAINT: computes a size constraint for  $\mathbf{H}_t$ ,
3. FILTERBYRANK: ranks and removes  $k$ -mers from the union until it fits the size constraint.

At the root, the library is flattened to a 1D array (PROJECTTABLE). We first describe our two main heuristics and then discuss details of other named functions in Algorithm 1.

#### Algorithm 1 KRANK algorithm. Functions are discussed in the text.

```

1: procedure BUILDLIBRARY( $t$ )
2:    $\mathbf{H}_t \leftarrow$  an empty table with  $2^{2h}$  rows
3:   if  $t \in \mathcal{S}$  then
4:      $\mathcal{K}_t \leftarrow$  EXTRACTKMERS( $t$ )
5:      $\mathbf{H}_t \leftarrow$  FILLTABLE( $\mathcal{K}_t$ )
6:     return  $\mathbf{H}_t$ 
7:   for all  $t' \in \text{child}(t)$  do
8:      $\mathbf{H}_{t'} \leftarrow$  BUILDLIBRARY( $t'$ )
9:   for all  $t' \in \text{child}(t)$  do
10:     $\mathbf{H}_t \leftarrow$  UNIONTABLES( $\mathbf{H}_{t'}, \mathbf{H}_t$ )
11:   $M(t) \leftarrow$  ADAPTIVESIZECONSTRAINT( $t$ )
12:   $n \leftarrow |\mathbf{H}_t| - M(t)$ 
13:  if  $n > 0$  then
14:    FILTERBYRANK( $\mathbf{H}_t, n$ )
15:  if  $t = \text{root}$  then
16:    PROJECTTABLE( $\mathbf{H}_t$ )
17:  return  $\mathbf{H}_t$ 

```

#### ADAPTIVESIZECONSTRAINT

Let  $M(t) \leq M$  be the upper bound (i.e., the size constraint) for the number of  $k$ -mers kept for each table at each node (i.e., ADAPTIVESIZECONSTRAINT in Algorithm 1). Let  $r: \mathcal{T} \rightarrow (0, 1]$  assign a portion of the full budget  $M$  to each node  $t$ ; any function  $r$  is valid as long as it takes the value 1 at the root (i.e.,  $r(\text{root})=1$ ) and increases as we move up the tree (i.e.,  $r(t) \geq r(t')$  if  $t' \in \text{child}(t)$ ). We

evaluated two options for  $r(t)$ :

$$r(t) = \frac{\sum_{s \in \mathcal{S}_t} |\mathcal{K}_s|}{\sum_{s \in \mathcal{S}} |\mathcal{K}_s|} \quad (\text{total } k\text{-mer count constraint}), \text{ or} \quad (7)$$

$$r(t) = \frac{|\mathcal{S}_t|}{|\mathcal{S}|} \quad (\text{number of species constraint}).$$

The two options for  $r(t)$  would be identical if all species have equal-sized genomes and each species is sampled the same number of times ( $\mathcal{K}_s$  can be the union over multiple genomes). Their main difference is that the first option allows highly sampled species (or very large genomes) to contribute more, whereas the latter does not. With either definition, we then set

$$M(t) = \begin{cases} |\mathcal{K}_t| & t \in \mathcal{S}, \\ \sqrt{r(t)}M & \text{otherwise} \end{cases} \quad (8)$$

noting that  $M(t) \leq \sum_{t' \in \text{child}(t)} M(t')$  due to the concavity of the square root. The concavity matters because it ensures that the size constraint becomes more restrictive as we move up the tree. In this case, the allowed budget increases faster than proportionally until  $r(t) = 1/\sqrt{2} \approx 0.7$  and slower than proportionally after. Thus, at lower ranks, many of the decisions for removing  $k$ -mers are left to its ancestors; note that a  $k$ -mer not removed at a taxon can still be removed further up. Close to the leaves, very few  $k$ -mers would be removed (e.g., a single species out of  $10^4$  still is assigned 1% of the total budget when visiting this node), whereas close to the root, taxon budgets get closer to the full budget (e.g., a phyla with two-thirds of the species is assigned 82% of the budget). Also, going up the tree, the  $k$ -mer budget assigned to each taxon increases compared to each child, but it never exceeds the total budget assigned to its children. Clearly,  $M(\text{root})$  matches the full budget allowed.

#### FILTERBYRANK

For each hash taxon  $t$ , the FILTERBYRANK function of Algorithm 1 removes  $k$ -mers when the size of the table  $\mathbf{H}_t$  exceeds  $M(t)$ . First, each row is pruned to  $b$   $k$ -mers according to the strategy described below. If the table size is still greater than  $M(t)$ , we continue removing  $k$ -mers, one at a time from each row in the descending order of their sizes, and iterating until the size constraint is satisfied.

For a  $k$ -mer set  $\mathcal{K}$ , a taxonomy  $\mathcal{T}$ , and its corresponding species set  $\mathcal{S}$ , we define  $R, R': \mathcal{K} \times \mathcal{T} \rightarrow [0, |\mathcal{S}|]$  as

$$R(x, t) = |\{y: x \in \mathcal{K}_y, y \in \mathcal{S}_t\}| \quad (9)$$

and

$$R^*(x, t) = \sum_{t' \in \text{child}(t)} \mathbb{1}\{|R(x, t')| > 0\}, \quad (10)$$

both illustrated in Supplemental Figure S2. We express each ranking strategy using either function and removing low-ranked  $k$ -mers accordingly. Discriminative and common  $k$ -mer strategies simply rank  $k$ -mers inversely and proportionally (respectively) to  $R$  or  $R'$ . In the case of ties, which can happen frequently, especially near leaves, we simply break them randomly.

The goal of covering all species can be imposed as a complex set covering problem. Instead, in the interest of scalability, we opt for a simple approach using weighted sums. We set

$$R^*(x, t) = \sum_{t' \in \text{child}(t)} w_{t'}(x) \cdot R(x, t') \quad (11)$$

where  $w_{t'}$  down-weights groups that are highly sampled among surviving  $k$ -mers (see Supplemental Fig. S2). A natural choice of weights is  $w_{t'} = 1/|\mathbf{H}_{t'}|$ . However, because  $|\mathbf{H}_{t'}|$  is independent of the  $k$ -mer or the row being considered, this approach will consistently

down-weight the children with the largest tables; in the extreme case, it can remove all  $k$ -mers from a child altogether. Thus, we define the weight of a  $k$ -mer  $x$  locally for its LSH table row by setting it to be inversely proportional to the *coverage* of each child taxon among  $k$ -mers of that row:

$$w_{t'}(x) = (|\{x' : x' \in \mathbf{H}_{t'}, \text{LSH}(x) = \text{LSH}(x')\}|)^{-1}.$$

As a result, children covered with fewer *surviving*  $k$ -mers get higher weights, leading the process at the parent node to remove them less often (see Supplemental Fig. S2 for a toy example). However, because decisions are made per row, and the assignment of  $k$ -mers to rows is random, this avoids consistently removing the same species.

### EXTRACTKMERS

In extracting  $k$ -mers from all input genomes, we use minimizers to perform an initial subsampling by choosing the  $k$ -mer whose encoding has a MurmurHash3 (<https://github.com/aappleby/smhasher>) value that is the smallest in a local window of size  $w = k + 3$ . We adopt the left/right encoding of  $k$ -mers introduced by Rachtman et al. (2021); this encoding allows fast calculation of HD using just four instructions (a pop-count, an XOR, an OR, and a shift). For  $k > 16$ , these encodings require 64 bits. However,  $h$  many  $k$ -mer positions are already given by the hash index and do not need to be saved as a part of the encoding. In KRANK, unlike CONSULT-II, we save memory by dropping these positions; this reduces the encoding size to a 32-bit number when  $k - h \leq 16$ , which is the case in our default configurations ( $k = 30$  and  $h = 14$  or  $k = 29$  and  $h = 13$ ). Because each  $k$ -mer encoding is only 32 bits, we save the encodings in the hash table; note that this is different from CONSULT-II, which saved pointers to a separate *encoding* array.

### Populating LSH tables (FILLTABLE, UNIONTABLES, and PROJECTTABLE)

For each leaf  $t$ , FILLTABLE computes LSH and encoding values of each  $k$ -mer  $x \in \mathcal{K}_t$ . Based on  $h$  positions, it stores  $k$ -mer's compact encoding in the row indexed by the LSH value. We keep a counter initialized to 1 associated with each  $k$ -mer to enable computing later. For nonleaf taxa, UNIONTABLES combines hash table rows by taking the union of  $k$ -mers. Luckily, the union operation can be defined on a per-row basis as the same  $k$ -mers share LSH value; independent union operations are performed in parallel across rows. Union also updates  $R(x, t)$ . Furthermore, KRANK uses some of the available threads to perform postorder traversal in parallel.

Below root, KRANK saves rows of the hash tables using vectors that allow dynamic size changes during construction, switching to static structures only at the end. At the root, each row gets pruned to  $b$   $k$ -mers, again using the ranking. This allows converting the table to a memory-efficient flattened static 1D array of size  $2^{2h}b$ .

### Batching

Because we need to merge multiple tables at internal nodes, the required memory may exceed what is available. To tackle this issue, we used a batching approach. Except for the ADAPTIVE SIZE CONSTRAINT, all components of our algorithm can be performed independently on a per-row basis. Therefore, we divide rows of the tables into batches and build them separately and in parallel. KRANK has a parameter to set the number of batches, and each batch's size constraint is simply updated as  $M(t)$  divided by the number of batches.

## Data sets

### WoL-v1 data set (read classification)

We demonstrate the pros and cons of selection strategies using the microbial WoL-v1 data set (Zhu et al. 2019) composed of 10,575 genomes (accessible at <https://biocore.github.io/wol/download>). We excluded 100 archaeal genomes used by Rachtman et al. (2021) as queries and five genomes with IDs missing from NCBI from the reference set, leaving us with 569 archaeal and 9901 bacterial references. We chose 756 queries in total; these include 666 bacterial genomes added to RefSeq after WoL-v1 was constructed, 10 randomly selected bacteria from the reference set, and all 80 of 100 archaeal queries from Rachtman et al. (2021) with Jaccard index above 0 to some reference genome. The queries span a range of distances to the closest reference genome and are binned into novelty groups as such (Supplemental Fig. S1A). The query set covers 55 phyla and 396 genera, with some phyla highly represented (e.g., 253 from *Pseudomonadota*, 91 from *Bacillota*) but including several rare phyla (Supplemental Fig. S1B). We generated 150 bp reads using ART (Huang et al. 2012) at high coverage with the default Illumina error profile and then subsampled down to 66,667 reads for each query ( $\approx 2 \times$  coverage for these genomes). We use the NCBI taxonomy provided with the WoL-v1 data set as our reference taxonomy.

We used this data set to explore the various heuristics of KRANK and to also compare the final KRANK versus other methods. We examined the same query set in both analyses. We compared KRANK against CONSULT-II v0.4.0 (Şapcı et al. 2024), where  $k$ -mer selection is arbitrary, Kraken 2 (with default settings), and CLARK (Ounit et al. 2015) (with default settings,  $k = 31$ ; species rank), which uses discriminative  $k$ -mers. For KRANK, the classification algorithm is identical to CONSULT-II v0.5.0. We tested KRANK with two memory levels: KRANK-hs (high-sensitivity) sets  $k = 30$ ,  $h = 14$ ,  $b = 16$ , and uses 51.2 GB of memory, whereas KRANK-lw (lightweight) sets  $k = 29$ ,  $h = 13$ ,  $b = 16$ , and 12.8 GB of memory. For details, see the exact commands and tool versions in the Supplemental Information.

### Profiling on the CAMI-I high-complexity data set

This data set contains five different high-complexity samples, each of size 75 Gbp, which are simulated by mimicking the abundance distribution of the underlying microbial communities (Szczyrba et al. 2017). We tested the Bracken (Lu et al. 2017) extension of Kraken 2, which combines its results with Bayesian priors for better profiling. We used the same custom libraries built with the WoL-v1 data set, with default parameters as described in the WoL-v1 data set (read classification). Results for CONSULT-II, CLARK, and Bracken are retrieved from our prior work and are available at GitHub (<https://github.com/bo1929/shared.CONSUULT-II>).

### Profiling on CAMI-II marine and strain-madness data sets

We used the ten-sample marine data set (5 Gbp each, "marmg" available at <https://frl.publisso.de/data/frl:6425521/marine/>) and the 100-sample (2 Gbp each, "strmg" available at <https://frl.publisso.de/data/frl:6425521/strain/>) strain-madness data set, which are the two main data sets focused on abundance profiling in the CAMI-II challenge. For all tools except CONSULT-II and KRANK, we used CAMI-II results submitted to the challenge available at GitHub ([https://github.com/CAMI-challenge/second\\_challenge\\_evaluation](https://github.com/CAMI-challenge/second_challenge_evaluation)). We exclude versions of the methods that were submitted to the first challenge. CONSULT-II results were again retrieved from GitHub (<https://github.com/bo1929/shared.CONSUULT-II>). We built and used a KRANK library in high-

sensitivity setting ( $k=30$ ,  $h=14$ ,  $b=16$ ) using 72,766 genomes selected from the NCBI RefSeq snapshot as of January 8, 2019 (retrieved from [https://openstack.cibitec.uni-bielefeld.de:8080/swift/v1/CAMI\\_2\\_DATABASES/](https://openstack.cibitec.uni-bielefeld.de:8080/swift/v1/CAMI_2_DATABASES/)) by allowing each species to contribute at most 500 genomes (otherwise randomly selected among available genomes).

### Evaluation metrics

We describe evaluation metrics, referring the reader to the Supplemental Methods section “Evaluation metrics and tools used” for further commands and details.

### Read classification

We evaluated the read classification with respect to the NCBI taxonomy. Each prediction is separately evaluated across each taxonomic rank. When the reference set had at least one genome matching the query taxon at a particular rank, we call it a *positive*: *TP* if it is the correct taxon, *FP* if it is the incorrect taxon, and *FN* if the read is not classified at that rank. Similarly, at a particular rank, when the reference set lacks a genome from the query taxon, we call it a *negative*: *TN* if the read is not classified at that rank, *FP* if it is classified, which would necessarily be false. We ignored queries at ranks where the true taxonomic ID given by NCBI is a missing rank.

We use custom scripts to process the specific output format of each tool and compare it with the taxonomy used to build reference libraries. These scripts are available on GitHub (<https://github.com/bo1929/shared.KRANK>). We also provide distances of each query genome to the closest reference genome,  $d^*$  as estimated by Mash (Ondov et al. 2016), taxonomy files, and the outputs of evaluation scripts.

### Taxonomic profiling

To measure taxonomic profiling performances on both CAMI challenges, we reported the same metrics as those emphasized in original publications. We used OPAL (version 1.0.12) to compute the metrics with `-n` option (which renormalizes after removing the unclassified portion).

On CAMI-I, we used the main two metrics singled out by the open-community profiling assessment tool (OPAL, available on GitHub: <https://github.com/CAMI-challenge/OPAL>) (Meyer et al. 2019): the Bray–Curtis dissimilarity between the estimated profile and the true profile and Shannon’s equitability as a measure of alpha diversity. All the results are averaged across all five samples. On the CAMI-II Challenge, we followed the original paper by Meyer et al. (2022) and compared methods based on L1 norm error and weighted UniFrac error. Both L1 norm error and weighted UniFrac error measure the accuracy of the relative abundance estimations of taxa in a sample. The UniFrac metric is computed across the entire taxonomic tree, whereas the L1 norm is computed at each taxonomic rank. The weighted UniFrac has a maximum of 16 on this data set, so we report 16—UniFrac; similarly, the L1 norm can never be more than 2, so we report 2—L1. We report averages across 10 samples and 100 samples, respectively, for marine and strain-madness data sets.

### Data access

Simulated query sequences for read classification can be found at <https://ter-trees.ucsd.edu/data/krank/KRANK-queries.tar.gz>. Our results, metadata files and scripts used to create figures, and custom scripts used to evaluate read classification results are available at GitHub (<https://github.com/bo1929/shared.KRANK>). KRANK’s

software can be found at GitHub (<https://github.com/bo1929/KRANK>). We also share both the software and custom auxiliary scripts as Supplemental Code. The list of KRANK libraries used across our experiments can be found at <https://ter-trees.ucsd.edu/data/krank/>, a tutorial for running KRANK, and descriptions of these libraries are available at GitHub.

### Competing interest statement

The authors declare no competing interests.

### Acknowledgments

We acknowledge the San Diego Supercomputer Center (SDSC) at University of California, San Diego for providing HPC resources that have contributed to the research results reported within this paper. The funding for this work was provided by the National Institute of Health (NIH) grant number R35GM142725 and by a research grant by the Minderoo Foundation to S.M.

*Author contributions:* A.O.B.Ş. and S.M. conceived the study, designed the experiments, and wrote the manuscript. A.O.B.Ş. wrote the KRANK software and conducted the experiments.

### References

- Asnicar F, Thomas AM, Beghini F, Mengoni C, Manara S, Manghi P, Zhu Q, Bolzan M, Cumbo F, May U, et al. 2020. Precise phylogenetic analysis of microbial isolates and genomes from metagenomes using PhyloPhlAn 3.0. *Nat Commun* **11**: 2500. doi:10.1038/s41467-020-16366-7
- Balaban M, Jiang Y, Zhu Q, McDonald D, Knight R, Mirarab S. 2024. Generation of accurate, expandable phylogenomic trees with uDance. *Nat Biotechnol* **42**: 768–777. doi:10.1038/s41587-023-01868-8
- Har-Peled S, Indyk P, Motwani R. 2012. Approximate nearest neighbors: towards removing the curse of dimensionality. *Theory Comput* **8**: 321–350. doi:10.4086/toc.2012.v008a014
- Huang W, Li L, Myers JR, Marth GT. 2012. ART: a next-generation sequencing read simulator. *Bioinformatics* **28**: 593–594. doi:10.1093/bioinformatics/btr708
- Lee D, Karchin R, Beer MA. 2011. Discriminative prediction of mammalian enhancers from DNA sequence. *Genome Res* **21**: 2167–2180. doi:10.1101/gr.121905.111
- Liang Q, Bible PW, Liu Y, Zou B, Wei L. 2020. DeepMicrobes: taxonomic classification for metagenomics with deep learning. *NAR Genom Bioinform* **2**: lqaa009. doi:10.1093/nargab/lqaa009
- Lu J, Breitwieser FP, Thielen P, Salzberg SL. 2017. Bracken: estimating species abundance in metagenomics data. *PeerJ Comput Sci* **3**: e104. doi:10.7717/peerj-cs.104
- McDonald D, Jiang Y, Balaban M, Cantrell K, Zhu Q, Gonzalez A, Morton JT, Nicolaou G, Parks DH, Karst SM, et al. 2024. Greengenes2 unifies microbial data in a single reference tree. *Nat Biotechnol* **42**: 715–718. doi:10.1038/s41587-023-01845-1
- Meyer F, Bremges A, Belmann P, Janssen S, McHardy AC, Koslicki D. 2019. Assessing taxonomic metagenome profilers with OPAL. *Genome Biol* **20**: 51. doi:10.1186/s13059-019-1646-y
- Meyer F, Fritz A, Deng ZL, Koslicki D, Lesker TR, Gurevich A, Robertson G, Alser M, Antipov D, Beghini F, et al. 2022. Critical assessment of metagenome interpretation: the second round of challenges. *Nat Methods* **19**: 429–440. doi:10.1038/s41592-022-01431-4
- Nasko DJ, Koren S, Phillippy AM, Treangen TJ. 2018. RefSeq database growth influences the accuracy of  $k$ -mer-based lowest common ancestor species identification. *Genome Biol* **19**: 165. doi:10.1186/s13059-018-1554-6
- Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, Phillippy AM. 2016. Mash: fast genome and metagenome distance estimation using MinHash. *Genome Biol* **17**: 132. doi:10.1186/s13059-016-0997-x
- Ounit R, Lonardi S. 2016. Higher classification sensitivity of short metagenomic reads with CLARK-S. *Bioinformatics* **32**: 3823–3825. doi:10.1093/bioinformatics/btw542
- Ounit R, Wanamaker S, Close TJ, Lonardi S. 2015. CLARK: fast and accurate classification of metagenomic and genomic sequences using discriminative  $k$ -mers. *BMC Genomics* **16**: 236. doi:10.1186/s12864-015-1419-2
- Pachiadaki MG, Brown JM, Brown J, Bezuidt O, Berube PM, Biller SJ, Poulton NJ, Burkart MD, La Clair JJ, Chisholm SW, et al. 2019. Charting the

- complexity of the marine microbiome through single-cell genomics. *Cell* **179**: 1623–1635.e11. doi:10.1016/j.cell.2019.11.017
- Parks DH, Chuvochina M, Waite DW, Rinke C, Skarshewski A, Chaumeil PA, Hugenholtz P. 2018. A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat Biotechnol* **36**: 996–1004. doi:10.1038/nbt.4229
- Rachtman E, Balaban M, Bafna V, Mirarab S. 2020. The impact of contaminants on the accuracy of genome skimming and the effectiveness of exclusion read filters. *Mol Ecol Resour* **20**: 649–661. doi:10.1111/1755-0998.13135
- Rachtman E, Bafna V, Mirarab S. 2021. CONSULT: accurate contamination removal using locality-sensitive hashing. *NAR Genom Bioinform* **3**: lqab071. doi:10.1093/nargab/lqab071
- Roberts M, Hayes W, Hunt BR, Mount SM, Yorke JA. 2004. Reducing storage requirements for biological sequence comparison. *Bioinformatics* **20**: 3363–3369. doi:10.1093/bioinformatics/bth408
- Şapcı AOB, Rachtman E, Mirarab S. 2024. CONSULT-II: accurate taxonomic identification and profiling using locality-sensitive hashing. *Bioinformatics* **40**: btae150. doi:10.1093/bioinformatics/btae150
- Sczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Dröge J, Gregor I, Majda S, Fiedler J, Dahms E, et al. 2017. Critical assessment of metagenome interpretation—a benchmark of metagenomics software. *Nat Methods* **14**: 1063–1071. doi:10.1038/nmeth.4458
- von Meijenfeldt FAB, Arkhipova K, Cambuy DD, Coutinho FH, Dutilh BE. 2019. Robust taxonomic classification of uncharted microbial sequences and bins with CAT and BAT. *Genome Biol* **20**: 217. doi:10.1186/s13059-019-1817-x
- Wood DE, Lu J, Langmead B. 2019. Improved metagenomic analysis with Kraken 2. *Genome Biol* **20**: 257. doi:10.1186/s13059-019-1891-0
- Zheng H, Marçais G, Kingsford C. 2023. Creating and using minimizer sketches in computational genomics. *J Comput Biol* **30**: 1251–1276. doi:10.1089/cmb.2023.0094
- Zhu Q, Mai U, Pfeiffer W, Janssen S, Asnicar F, Sanders JG, Belda-Ferre P, Al-Ghalith GA, Kopylova E, McDonald D, et al. 2019. Phylogenomics of 10,575 genomes reveals evolutionary proximity between domains Bacteria and Archaea. *Nat Commun* **10**: 5477. doi:10.1038/s41467-019-13443-4

Received March 15, 2024; accepted in revised form August 6, 2024.