

Lightweight compositional analysis of metagenomes with sourmash gather

This manuscript ([permalink](#)) was automatically generated from [dib-lab/2020-paper-sourmash-gather@4e66936](#) on October 29, 2021.



Authors

- **Luiz Irber**

 [0000-0003-4371-9659](#) ·  [luizirber](#) ·  [luizirber](#)

Graduate Group in Computer Science, UC Davis; Department of Population Health and Reproduction, UC Davis ·
Funded by Grant XXXXXXXX

- **C. Titus Brown**

 [0000-0001-6001-2677](#) ·  [ctb](#)

Department of Population Health and Reproduction, UC Davis

Abstract

The assignment of genomes and taxonomy to metagenome data underlies many microbiome studies. Here we describe two algorithms for compositional analysis of metagenome sequencing data. We first develop a sketching technique, *Scaled MinHash*, that supports Jaccard containment estimation. We implement *Scaled MinHash* in the sourmash software and demonstrate large-scale containment searches of metagenomes using all 700,000 currently available microbial reference genomes. We next frame shotgun metagenome compositional analysis in terms of min-set-cover, i.e. as the problem of finding a minimum collection of reference genomes that “cover” the known portion of a metagenome. We implement a greedy approximate solution using *Scaled MinHash* sketches, and evaluate its accuracy in taxonomic assignment using a CAMI community benchmark. Finally, we show that the minimum set cover can be used to for read mapping. sourmash is available as open source under the BSD 3-Clause license at github.com/dib-lab/sourmash/.

Introduction

Shotgun metagenomics samples the DNA sequence content of microbial communities.

Compositional analysis of shotgun metagenome samples has the goal of identifying what reference genomes to use for functional and taxonomic interpretation of metagenome content.

The substantial increase in the number of available reference genomes presents a significant practical obstacle to comprehensive compositional analyses.

Here, we describe a lightweight approach to compositional analysis of shotgun metagenome data. Our approach tackles the selection of appropriate reference genomes and provides a computationally efficient method for taxonomic classification of metagenome data.

We first define *Scaled MinHash*, an extension of MinHash sketching that supports lightweight containment estimation for metagenome datasets using k-mers. We implement *Scaled MinHash* in a Python and Rust package, `sourmash`, and show that it is competitive in accuracy with other containment estimation approaches.

We next frame reference-based metagenome content analysis as a min-set-cov problem, where we determine the *minimum* number of genomes from a reference database needed to cover the identifiable genomic content from a metagenome. We implement a best-polynomial-time greedy approximation to the min-set-cov problem using *Scaled MinHash* in `sourmash`, provides an interactive decomposition of metagenomes into genome matches.

To evaluate the accuracy of our min-set-cov procedure, we implement a simple taxonomic classification approach in which we use the taxonomy of the genomes from the set cover to define the taxonomy of the metagenome content. We show that this permits precise and lightweight classification of metagenome content across all taxonomic levels.

Finally, we show that the minimum set covers for several metagenomes are only a small subset of genomes even when using very large and redundant databases, and demonstrate that this subset can be used to map the metagenome reads in concordance with the estimates from *Scaled MinHash*. Thus, *Scaled MinHash* combined with min-set-cov provides a lightweight, accurate, and useful way to estimate the composition of metagenomes using a large reference database.

Results

Scaled MinHash sketches support accurate containment operations

We define the Scaled MinHash on an input domain of k -mers, W , as follows:

$$\mathbf{SCALED}_s(W) = \{ w \leq \frac{H}{s} \mid \forall w \in W \}$$

where H is the largest possible value in the domain of $h(x)$ and $\frac{H}{s}$ is the value in the Scaled MinHash.

The Scaled MinHash is a mix of MinHash and ModHash [1]. It keeps the selection of the smallest elements from MinHash, while using the dynamic size from ModHash to allow containment estimation. However, instead of taking $0 \bmod m$ elements like $\mathbf{MOD}_m(W)$, a Scaled MinHash uses the parameter s to select a subset of W .

Scaled MinHash supports containment estimation with high accuracy and low bias. (Analytic work from David HERE.)

- approximation formula (eqn 13 from overleaf)
- for queries into large sets (large $|A|$), bias factor is low.
- refer to appendix for derivation.

Given a uniform hash function h and $s = m$, the cardinalities of $\mathbf{SCALED}_s(W)$ and $\mathbf{MOD}_m(W)$ converge for large $|W|$. The main difference is the range of possible values in the hash space, since the Scaled MinHash range is contiguous and the ModHash range is not. This permits a variety of convenient operations on the sketches, including iterative downsampling of Scaled MinHash sketches as well as conversion to MinHash sketches.

A Scaled MinHash implementation accurately estimates containment between sets of different sizes

We compare the *Scaled MinHash* method to CMash (*Containment MinHash*) [2] and Mash Screen (*Containment Score*) [3] for containment queries in a synthetic mock metagenomic bacterial and archaeal community where the reference genomes are largely known [4]. This data set has been used in several methods evaluations [3].



Figure 1: Letter-value plot [hofmann letter-value 2017?] of the differences from containment estimate to ground truth (exact). Each method is evaluated for $k = \{21, 31, 51\}$, except for Mash with $k = 51$, which is unsupported.

Figure 1 shows results with low-coverage and contaminant genomes (as described in [7] and [8]) removed from the database. All methods are within 1% of the exact containment on average (Figure 1), with CMash consistently underestimating the containment for large k and overestimating for small k . Mash Screen with $n = 10000$ has the smallest difference to ground truth for $k = \{21, 31\}$, followed by smol with scaled=1000 and Mash Screen with $n = 1000$.

CTB todo:

- use sourmash, not smol

CTB questions:

- should we add sketch sizes in here more explicitly? e.g. number of hashes kept?
- compares well with others
- How much is missed figure; Poisson calculations? => appendix?
- should we make comment here about sequencing errors?

We can use Scaled Min-Hash to construct a minimum set cover for metagenomes

We next ask: what is the smallest collection of genomes in a database that contains all of the known k-mers in a metagenome? Formally, for a given metagenome M and a reference database D , what is the minimum collection of genomes in D which contain all of the k-mers in the intersection of D and M ? That is, we wish to find the smallest set $\{G_n\}$ of genomes in D such that

$$k(M) \cap k(D) = \bigcup_n \{k(M) \cap k(G_n)\}$$

This is the *minimum set covering* problem, for which there is a polynomial-time approximation (cite). (Provide algorithm here.)

This greedy algorithm works by iteratively subtracting k-mers belonging to the genome that has the highest containment count from the metagenome (ref alg above). This results in a progressive classification of the known k-mers in the metagenome to specific genomes, in rank order of number of contained hashes. Note that in cases where equivalent matches are available at a particular rank, the match is chosen at random.

In Figure 2, we show the results of this iterative decomposition of the mock metagenome from [4], into constituent genome matches. The high rank (early) matches reflect large and/or mostly-covered genomes with high containment, while later matches reflect smaller genomes, lower-covered genomes, and/or genomes with substantial overlap with earlier matches. Where there are overlaps between genomes, shared common k-mers are “claimed” by higher rank matches and only k-mer content specific to the later genome is used to identify the lower rank matches. For example, genomes from two strains of *Shewanella baltica* present in the mock metagenome in Figure 2 have an approximately 50% overlap in k-mer content, and these shared k-mers are claimed by *Shewanella baltica* OS223 (compare *Shewanella baltica* OS223, rank 8, with *Shewanella baltica* OS185, rank 33; the difference between the red circles and green triangles for *S. baltica* OS185 represents the k-mers claimed by *S. baltica* OS223). CTB: maybe indicate or highlight these genomes in the figure?

For this mock metagenome, 205m (54.8%) of 375m k-mers were found in GenBank. The remaining 169m (45.2%) k-mers had no matches, and represent either k-mers introduced by sequencing errors or unknown k-mers from real community members.



Figure 2: K-mer decomposition of a metagenome into constituent genomes. A rank ordering for the first 36 genomes from the minimum set cover of the synthetic metagenome from [4], calculated using 700,000 GenBank genomes. The Y axis is labeled with the NCBI-designed name of the genome. In the left plot, the X axis represents the estimated number of k-mers shared between each genome and the metagenome. The red circles indicate the number of remaining k-mers at that rank, while the green triangle symbols indicate the total number of k-mers, including those already assigned at previous ranks. In the right plot, the X axis represents the estimated k-mer coverage of that genome. The red circles indicate the coverage with the remaining k-mers at that rank, while the green triangle symbols indicate total coverage with all k-mers in the metagenome, including those already assigned at previous ranks.

Minimum metagenome covers can accurately estimate taxonomic composition

We evaluated the accuracy of min-set-cov for metagenome decomposition by using benchmarks from the Critical Assessment of Metagenome Interpretation (CAMI) [9], a community-driven initiative for reproducibly benchmarking metagenomic methods. We used the mouse gut metagenome dataset [10], in which a simulated mouse gut metagenome (MGM) was derived from 791 bacterial and archaeal genomes, representing 8 phyla, 18 classes, 26 orders, 50 families, 157 genera, and 549 species. 64 samples were generated with CAMISIM, with 91.8 genomes present on each sample on average. Each sample is 5 GB in size, and both short-read (Illumina) and long-read (PacBio) simulated sequencing data is available. (CTB: check citations / content of latest actual CAMI pub, <https://www.biorxiv.org/content/10.1101/2021.07.12.451567v1>)

Since min-set-cov yields only a collection of genomes rather than a species list, we generated a taxonomic profile for a given metagenome cover through the following procedure. For each genome match, we used the species designation in the NCBI taxonomy for that genome. Then, we calculated the fraction of the genome remaining in the metagenome after k-mers belonging to higher-rank genomes have been removed (red circles in Figure 2 (a)). We used this fraction to weight the contribution of the genome's species designation towards the metagenome taxonomy. This procedure produces an estimate of that species' taxonomic contribution to the metagenome, normalized by the genome size.

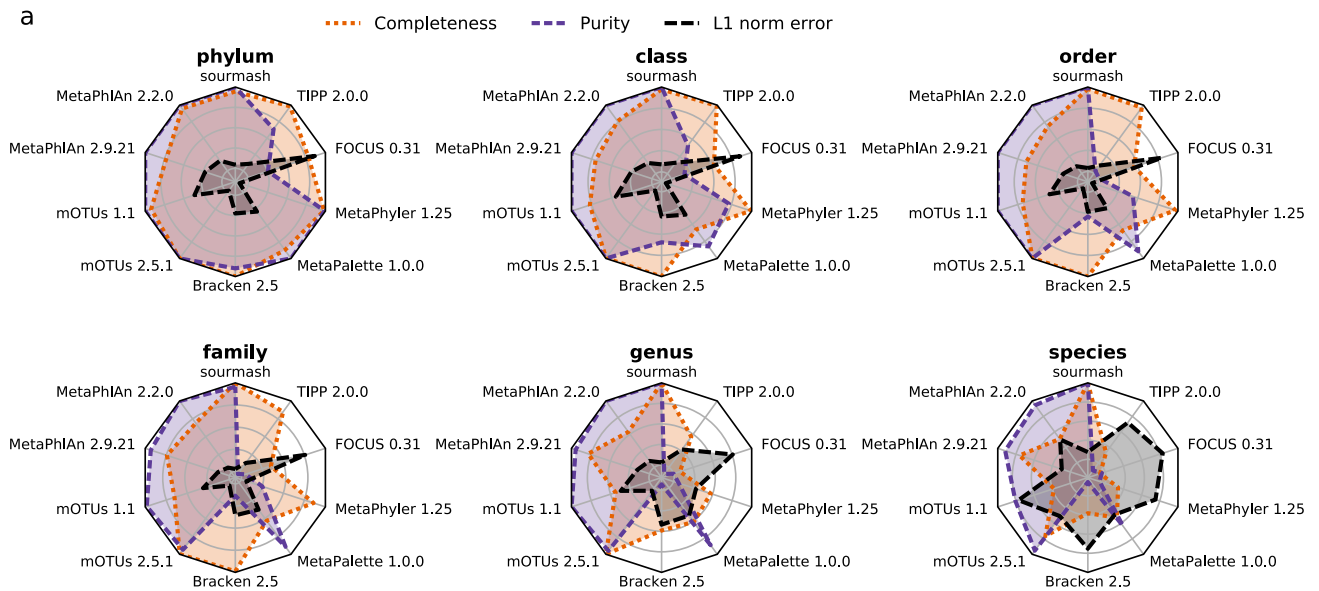


Figure 3: Comparison per taxonomic rank of methods in terms of completeness, purity (1% filtered), and L1 norm.

C	Completeness	Purity (1% filtered)	L1 norm error	Sum of scores
1st	sourmash (247)	sourmash (179)	mOTUs 2.5.1 (789)	sourmash (1262)
2nd	mOTUs 2.5.1 (416)	MetaPhlAn 2.2.0 (241)	sourmash (836)	mOTUs 2.5.1 (1887)
3rd	Bracken 2.5 (1008)	mOTUs 1.1 (631)	MetaPhlAn 2.9.21 (1401)	MetaPhlAn 2.2.0 (3527)
4th	MetaPhyler 1.25 (1298)	mOTUs 2.5.1 (682)	MetaPhlAn 2.2.0 (1497)	MetaPhlAn 2.9.21 (4349)
5th	TIPP 2.0.0 (1424)	MetaPhlAn 2.9.21 (789)	MetaPhyler 1.25 (1586)	MetaPhyler 1.25 (5148)
6th	MetaPhlAn 2.2.0 (1789)	MetaPalette 1.0.0 (1182)	mOTUs 1.1 (2317)	mOTUs 1.1 (5253)
7th	MetaPhlAn 2.9.21 (2159)	MetaPhyler 1.25 (2264)	TIPP 2.0.0 (2361)	MetaPalette 1.0.0 (5989)
8th	mOTUs 1.1 (2305)	Bracken 2.5 (2881)	MetaPalette 1.0.0 (2390)	Bracken 2.5 (6574)
9th	MetaPalette 1.0.0 (2417)	TIPP 2.0.0 (3361)	Bracken 2.5 (2685)	TIPP 2.0.0 (7146)
10th	FOCUS 0.31 (3424)	FOCUS 0.31 (3764)	FOCUS 0.31 (3894)	FOCUS 0.31 (11082)

Figure 4: Methods rankings and scores obtained for the different metrics over all samples and taxonomic ranks. For score calculation, all metrics were weighted equally.

In Figures 3 and 4 we show an updated version of Figure 6 from [10] that includes our method, implemented in the `sourmash` software (CTB: what databases are used?). Here we compare 10 different methods for taxonomic profiling and their characteristics at each taxonomic rank. While previous methods show reduced completeness, the ratio of taxa correctly identified in the ground truth, below the genus level, `sourmash` can reach 88.7% completeness at the species level with the highest purity (the ratio of correctly predicted taxa over all predicted taxa) across all methods: 95.9% when filtering predictions below 1% abundance, and 97% for unfiltered results. `sourmash` also has the lowest L1-norm error (the sum of the absolute difference between the true and predicted abundances at a specific taxonomic rank), the highest number of true positives and the lowest number of false positives.

Minimum metagenome covers select small subsets of large databases

Table 1: metagenomes and min-set-cov.

data set	genomes >= 100k overlap	min-set-cov	% k-mers identified
zymo mock (SRR12324253)	405,839	19	0%
podar mock (SRR606249)	5800	74	0%
gut real (SRR5650070)	96,423	99	0%
oil well real (SRR1976948)	1235	135	0%

In Table 1, we show the results of running min-set-cov for four metagenomes against genbank - two mock communities (cite cite), one human gut microbiome data set from iHMP (cite), and an oil well sample (cite). Our implementation provides estimates for both the *total* number of genomes with substantial overlap to a query genome, and a *minimum list* of genomes that account for k-mers with overlap in the query metagenome (see Methods).

We find many genomes with large overlaps for each metagenome, due to the redundancy of the reference database. For example, the zymo mock contains a *Salmonella* genome, and there are over 200,000 *Salmonella* genomes that match to it in Genbank. Likewise, the iHMP dataset contains many XYZ. Since neither the podar mock nor the oil well community contain genomes from species with substantial representation in genbank, they yield many fewer total overlapping genomes.

However, regardless of the number of genome with overlap, the estimated *minimum* collection of genomes is always much smaller than the number of genomes with overlaps. In the cases where the k-mers in the metagenome are mostly identified, this is because of database redundancy: e.g. in the case of the zymo mock, the min-set-cov algorithm chooses only one *Salmonella* genome from the 200,000+ available. Conversely, in the case of the oil well sample, much of the sample is not identified, suggesting that the small size of the covering set is because much of the sample is not represented in the database.

CTB TODO: add % identified to table!

Minimum metagenome covers provide representative genomes for mapping

Mapping metagenome reads to representative genomes is an important step in many microbiome analysis pipelines, but mapping approaches struggle with large, redundant databases. One use case for a minimum metagenome cover is to select a small set of representative genomes to be used for mapping. We therefore developed a hybrid selection and mapping pipeline that uses the rank-ordered min-set-cov results to map reads to candidate genomes.

We first map all reads to all genomes in the minimum set cover, and then successively remove reads that map to higher rank genomes from lower rank genomes, and remap the remaining reads. That is, all reads mapped to the rank-1 genome in Figure 2 are removed from the rank-2 genome mapping, and all reads mapping to rank-1 and rank-2 genomes are removed from the rank-3 genome mapping. This produces results directly analogous to those presented in Figure 2, but for reads rather than k-mers (CTB: provide as Suppl Figure?). Importantly, in this process we only consider genomes identified

in the minimum set cover, because it is computationally intractable to map reads to the entire GenBank database. (CTB: check centrifuge?)

Figure 5 compares hash assignment rates and mapping rates for the four evaluation metagenomes in Table 1. Broadly speaking, we see that k-mer based estimates of metagenome composition align closely with the number of bases covered by mapped reads. This suggests that the k-mer based min-set-cov approach effectively selects reference genomes for metagenome read mapping.

For mock metagenomes (panels X and Y), there appears to be a close correspondence between mapping and hash assignment rates, while for actual metagenomes, there is more variation between mapping and hash assignments. Further work is needed to evaluate rates of variation across a larger number of metagenomes.

CTB: do this for all four metagenomes!

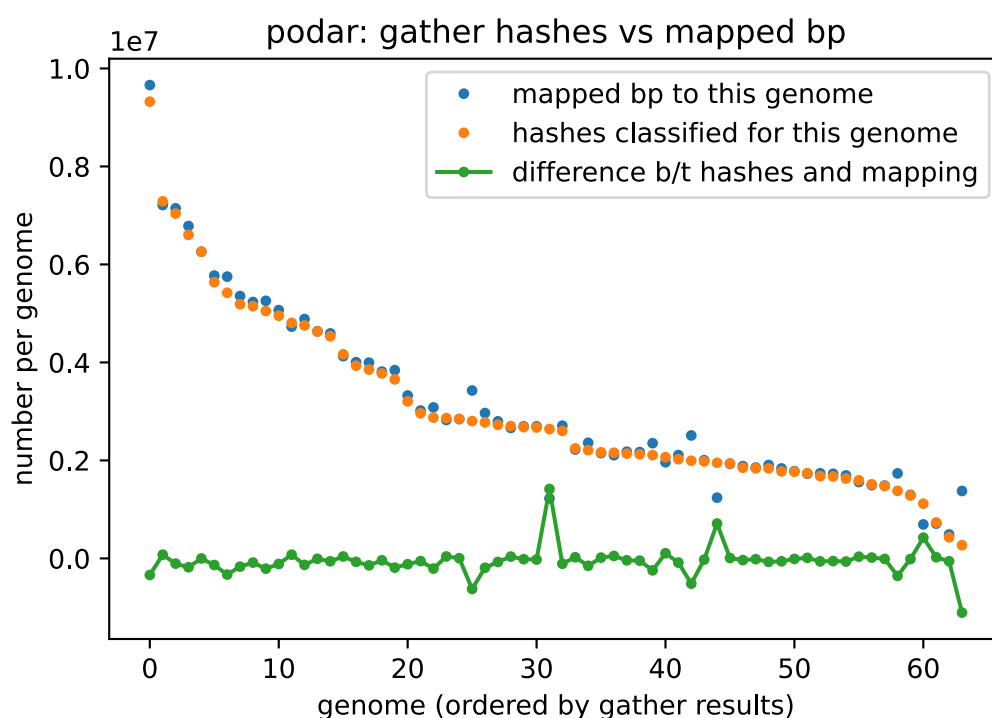


Figure 5: Hash-based decomposition of a metagenome into constituent genomes compares well to bases covered by read mapping. The reference genomes are rank ordered along the x axis based on the largest number of hashes from the metagenome specific to that genome; hence the number of hashes classified for each genome (orange dots) is monotonically decreasing. The y axis shows absolute number of estimated k-mers classified to this genome (orange) or total number of bases covered in the reference (blue); the numbers have not been rescaled. Decreases in mapping (green peaks) occur for genomes which are not exact matches to the genomes of the organisms used to build the mock community (cite sherine, mash screen).

Discussion

Scaled MinHash provides efficient containment queries for large data sets.

Scaled MinHash is an implementation of ModHash that uses the bottom hashing concept from MinHash to support containment operations. In brief, all elements in the set to be sketched are

hashed, and any hash values below a certain fixed boundary value are kept for the sketch. This fixed boundary value is determined by the desired accuracy for the sketch representation.

Intuitively, *Scaled MinHash* performs a density sampling at a rate of 1 k -mer per s distinct k -mers seen, where s is the size of the hash space divided by the boundary value used in creating the sketch. This is a type of lossy compression, with a fixed compression ratio of s : for values of s used here ($s \approx 1000$), data sets are reduced in size 1000-fold.

Unlike MinHash, *Scaled MinHash* supports containment analysis between sets of very different sizes, and here we demonstrate that it can be used efficiently and effectively for compositional analysis of shotgun metagenome data sets with k -mers. In particular, *Scaled MinHash* is competitive in accuracy with extant MinHash-based techniques for containment analysis, while also supporting Jaccard similarity. Footnote: We note that others have also applied the ModHash concept to genomic data; see, for example, Durbin's "modimizer" [11].

Scaled MinHash offers several conveniences over *MinHash*. No hash is ever removed from a *Scaled MinHash* sketch during construction; while this means that sketches grow proportionally to the number of distinct k -mers in the sampled data set, sketches *also* support many operations - including all of the operations used in this paper - without needing to revisit the original data set. This is in contrast to MinHash, which requires auxiliary data structures for many operations - most especially, containment operations (cite CMash and mash screen). Thus Scaled MinHash sketches serve as distributed compressed indices for the original content for a much broader range of operations than MinHash.

Because *Scaled MinHash* sketches collect all hash values below a fixed threshold, they also support streaming analysis of sketches: any operations that used a previously selected value can be cached and updated with newly arriving values. ModHash has similar properties, but this is not the case for MinHash, since after n values are selected any displacement caused by new data can invalidate previous calculations.

Scaled MinHash also directly supports the addition and subtraction of hash values from a sketch, allowing post-processing and filtering without revisiting the original data set. This includes unions and intersections. Although possible for MinHash, in practice this requires oversampling (using a larger n) to account for possibly having less than n values after filtering; this approach is taken by Finch, another MinHash sketching software for genomics [12].

When the multiplicity of hashes in the original data is retained, *Scaled MinHash* sketches can be filtered on abundance. This allows removing low-abundance values, as implemented in Finch [doi:10.21105/joss.00505]. Filtering values that only appear once was implemented in Mash by using a Bloom Filter and only adding values after they were seen once, with later versions also implementing an extra counter array to keep track of counts for each value in the MinHash. These operations can be done in *Scaled MinHash* without auxiliary data structures.

Another useful operation available on *Scaled MinHash* sketches is *downsampling*: the contiguous value range for Scaled MinHash sketches allows deriving MinHash sketches from *Scaled MinHash* sketches whenever the number of hashes in the *Scaled MinHash* sketch is equal to or greater than n , as long as the same hashing scheme is used. Likewise, MinHash sketches can be converted to *Scaled MinHash* sketches when the maximum hash value in the MinHash sketch is larger than s .

Finally, because *Scaled MinHash* sketches are simply collections of hashes, existing k -mer indexing approaches can be applied to the sketches to provide fast database search of these indices with both similarity and containment metrics; a number of indexing operations, including Sequence Bloom Trees and reverse indices, are provided in the sourmash software.

In exchange for these many conveniences, *Scaled MinHash* sketches have limited sensitivity for small data sets where the k-mer cardinality of the data set $\approx s$, and are only bounded in size by H/s (typically quite large, $\approx 2e16$). The limited sensitivity of sketches may affect the sensitivity of gene- and viral genome-sized queries, but at $s = 1000$ we see comparable accuracy and sketch size to MinHash for bacterial genome comparisons (Figure XXX, currently 1).

Minimum set covers can be used for accurate compositional analysis of metagenomes.

Many metagenome content analysis approaches use reference genomes to interpret metagenome content, but most such approaches rely on a curated list of non-redundant genomes from a much larger database (e.g. Humann3/biobakery, cite). Here, we seek the *minimum* set of reference genomes necessary to account for all k-mers shared between the metagenome and the reference database. We show that this can be resolved efficiently for real-world data sets using a greedy algorithm; using *Scaled MinHash*, we provide an approach that readily scales to 700,000 genomes on current hardware (performance in appendix). Moreover, this procedure reduces the number of genomes under consideration to ≈ 100 for several mock and real metagenomes.

The development of a small list of relevant genomes is particularly useful for large reference databases containing many redundant genomes; for example, in Table XXX, we show that for one particular mock community, we can select a minimum metagenome cover of 19 genomes for a metagenome that contains matches to over 400,000 genomes total.

This minimum metagenome cover can then be used as inputs for further analysis, including both taxonomic content analysis and mapping approaches. For taxonomic analyses, we find that this approach is competitive with other current approaches and has many additional conveniences (discussed in detail below). The comparison of hash-based estimation of containment to mapping results in Figure ?? shows that this approach is an accurate proxy for systematic mapping.

This min-set-cov approach for assigning genomes to metagenomes using k-mers differs substantially from extant k-mer and mapping-based approaches for identifying relevant genomes. LCA-based approaches such as Kraken label individual k-mers based on taxonomic lineages in a database, and then use the resulting database of annotated k-mers to assign taxonomy to reads. Mapping- and homology-based approaches such as Diamond or @@@ use read mapping to genomes or read alignment to gene sequences in order to assign taxonomy and function. These approaches typically focus on assigning *individual* k-mers or reads. In contrast, here we analyze the entire collection of reads/k-mers and assigns them *in aggregate* to the *best* genome match.

Our implementation of the min-set-cov algorithm in sourmash also readily supports custom reference databases as well as updating minimum set covers with the addition of new reference genomes. When updating set covers, the first stage of calculating overlaps can be done on only the new genomes (Column XYZ of Table ZZZ), while the actual calculation of the minimum set cover must be redone.

Our implementation of min-set-cov on top of *Scaled MinHash* means that there is a loss of resolution when choosing between very closely related genomes, because distinct hashes will not be chosen for them. However, other data structures can be easily used for min-set-cov: any data structure supporting both the *containment* $C(A, B) = \frac{|A \cap B|}{|A|}$ and *remove elements* operations can be used to implement the greedy approximation algorithm (ref algorithm in results section 1). For example, a simple *set* of the *k*-mer composition of the query supports element removal, and calculating containment can be done with regular set operations. Approximate membership query (AMQ)

sketches like the *Counting Quotient Filter* [[pandey_general-purpose 2017?](#)] can also be used, with the benefit of reduced storage and memory usage. Moreover, the collection of datasets can be implemented with any data structure that can do containment comparisons with the query data structure. It is important to have performant containment searches, since `gather` may run `FindBestContainment` many times.

The min-set-cov approach is reference-based, and hence is very dependent on the reference database. In particular, in many cases the exact reference strains present in the metagenome will not be present in the database. This manifests in two ways in Figure ???. First, there is a systematic mismatch between the hash content and the mapping content (green line), because mapping software is more permissive in the face of small variants than k-mer-based exact matching. Moreover, many of the lower rank genomes in the plot are from the same species but different *strains* as the higher ranked genomes, suggesting that strain-specific portions of the reference are being utilized for matching at lower ranks. In reality, there will usually be a different mixture of strains in the metagenome than in the reference database. Approaches such as spacegraphcats may help resolve this by adapting old references. [[cite?](#)].

CTB: can we guess at places where gather would break? One is equivalent containment/different genome sizes, e.g. virus/phage contained within other genomes.

Minimum metagenome covers support accurate and flexible taxonomic conversation

Once the min-set-cov approach has identified reference genomes, we can build a taxonomic classifier for metagenome content by simply reporting the taxonomies of the constituent genomes, aggregated at the relevant taxonomic level using an LCA approach. Our initial taxonomic benchmarking shows that this approach is competitive for all metrics across all taxonomic levels (Figure XYZ).

One convenient feature of this approach to taxonomic analysis is that new or changed taxonomies can be readily incorporated by assigning them directly to genome identifiers; the majority of the computational work is involved in finding the reference genomes, which can have assignments in different taxonomic frameworks. For example, sourmash already supports GTDB natively, and will also support the emerging LINS framework (cite, cite). sourmash can also readily incorporate updates to taxonomies, e.g. the frequent updates to the NCBI taxonomy, without requiring expensive reanalysis of the primary metagenome data or the min-set-cov computation. Interestingly, the framing of taxonomic classification as a minimum set cover problem may also avoid the loss of taxonomic resolution that affects k-mer- and read-based approaches (cite cite); this is because we apply LCA *after* reads and k-mers have been assigned to individual genomes, and choose entire *genomes* based on a greedy best-match-first approach.

Finally, as the underlying min-set-cov implementation supports custom databases, it is straightforward to support taxonomic analysis using custom databases and/or custom taxonomic assignments. (sourmash already supports this natively.)

Algorithm is simple, computational performance is great

The algorithms underlying both *Scaled MinHash* and the greedy min-set-cov approximation are simple to describe and straightforward to implement. This increases the likelihood of correct implementation, provides opportunities for independent optimization of data structures, and simplifies interoperability between different implementations.

In the sourmash software package, we provide a mature and optimized implementation that implements all of the operations above. sourmash performs well in practice and supports a wide variety of use cases (see performance in appendix, docs and tutorials at sourmash.rtfd.io, and installation instructions for pip and conda). The sourmash project also provides large scale databases for NCBI and GTDB taxonomies.

Limitations of our approach

(For *Scaled MinHash*, `gather`, and taxonomy. Move where? Conclusions?)

`gather` as implemented in `sourmash` has the same limitations as *Scaled MinHash* sketches, including reduced sensitivity to small genomes/sequences such as viruses. *Scaled MinHash* sketches don't preserve information about individual sequences, and short sequences using large scaled values have increasingly smaller chances of having any of its k -mers (represented as hashes) contained in the sketch. Because it favors the best containment, larger genomes are also more likely to be chosen first due to their sketches have more elements, and further improvements can take the size of the match in consideration too. Note that this is not necessarily the *similarity* $J(A, B)$ (which takes the size of both A and B), but a different calculation that normalizes the containment considering the size of the match.

`gather` is also a greedy algorithm, choosing the best containment match at each step. Situations where multiple matches are equally well contained or many datasets are very similar to each other can complicate this approach, and additional steps must be taken to disambiguate matches. The availability of abundance counts for each element in the *Scaled MinHash* is not well explored, since the process of *removing elements* from the query doesn't account for them (the element is removed even if the count is much higher than the count in the match). Both the multiple match as well as the abundance counts issues can benefit from existing solutions taken by other methods, like the *species score* (for disambiguation) and *Expectation-Maximization* (for abundance analysis) approaches from Centrifuge [[kim centrifuge 2016?](#)].

(From David Koslicki) Gotchas:

- Lack of sensitivity for small queries
- Potentially large sketch sizes

And a couple other that I've tentatively/mathematically observed:

- The variance of the estimate of $C(A,B) = |AB| / |A|$ appears to also depend on $|A|$, which was somewhat surprising
- The "fixed k-size" problem (which might be able to be overcome with the prefix-lookup data structure, if one sacrifices some accuracy)

We belieeeeeeeve

min set cov could be applied in many more circumstances - read based analysis, contig based analysis, maybe variant calling, etc.

Mention weighted cover cc David?

Note that here we are providing one approach / approximation (Scaled MinHash containment) with one shingling approach (k-mers) to tackle metagenome composition for mapping and taxonomy. The

min-set-cover approach could be used with exact containment, and/or with other shingling approaches.

CTB: discuss centrifuge, etc. Could this be implemented on top of that?

Conclusion

- scaled min hash is powerful, with well defined limitations.
- gather is awesome and convenient.
- taxonomy is awesome and overcomes limitations of many current approaches.
- sourmash is robust software that provides a practically usable implementation of these ideas.
- future directions...

Scaled MinHash sketches are simple to implement and analyze, with consistent guarantees for the range of values and subsetting properties when applied to datasets. Containment and similarity operations between *Scaled MinHash* sketches avoid the need to access the original data or more limited representations that only allow membership query, and serve as a proxy for large scale comparisons between hundreds or thousands of datasets.

Small genomes require low scaled values in order to properly estimate containment and similarity, and exact k -mer matching is brittle when considering evolutionarily-diverged organisms. While some of these problems can be overcome in future work, *Scaled MinHash* sketches can serve as a prefilter for more accurate and computationally expensive applications, allowing these methods to be used in larger scales by avoiding processing data that is unlikely to return usable results.

Scaled MinHash sketches are effective basic building blocks for creating a software ecosystem that allow practical applications, including taxonomic classification in metagenomes and large scale indexing and searching in public genomic databases.

Methods

Implementation of Scaled MinHash

We provide two implementations of Scaled MinHash, `smol` and `sourmash`. `smol` is a minimal implementation of *Scaled MinHash* developed to demonstrate the method; it does not include many required features for working with real biological data, but its smaller code base makes it a more readable and concise example of the method. `sourmash` [13] implements features and functionality needed for large scale analyses of real data.

Comparison between CMash, mash screen, and Scaled MinHash.

Experiments use $k = \{21, 31, 51\}$ (except for Mash, which only supports $k \leq 32$). For Mash and CMash they were run with $n = \{1000, 10000\}$ to evaluate the containment estimates when using larger sketches with sizes comparable to the Scaled MinHash sketches with *scaled* = 1000. The truth set is calculated using an exact k -mer counter implemented with a *HashSet* data structure in the Rust programming language [matsakis rust 2014?].

For *Mash Screen* the ratio of hashes matched by total hashes is used instead of the *Containment Score*, since the latter uses a k -mer survival process modeled as a Poisson process first introduced in [fan assembly 2015?] and later used in the *Mash distance* [ondov mash: 2016?] and *Containment score* [ondov mash 2019?] formulations.

MHBT

The *MinHash Bloom Tree (MHBT)* is a variation of the *Sequence Bloom Tree (SBT)* that uses Scaled MinHash sketches as leaf nodes instead of Bloom Filters as in the SBT. The search operation in SBTs is defined as a breadth-first search starting at the root of the tree, using a threshold of the original k -mers in the query to decide when to prune the search. MHBTs use a query Scaled MinHash sketch instead, but keep the same search approach. The threshold of a query Q approach introduced in [solomon fast 2016?] is equivalent to the containment

$$C(Q, S) = \frac{|Q \cap S|}{|S|}$$

described in [broder resemblance 1997?], where S is a Scaled MinHash sketch. For internal nodes n (which are Bloom Filters) the containment of the query Scaled MinHash sketch Q is

$$C(Q, n) = \frac{|\{h \in n \mid \forall h \in Q\}|}{|Q|}$$

as defined by [koslicki improving 2019?] for the *Containment MinHash to Bloom Filter* comparison.

MHBTs support both containment and similarity queries. For internal nodes the containment $C(Q, n)$ is used as an upper-bound of the similarity $J(Q, n)$:

[Math Processing Error]

since $|Q \cup n| \geq |Q|$. When a leaf node is reached then the similarity $J(Q, S)$ is calculated for the Scaled MinHash sketch S and declared a match if it is above the threshold t . Because the upper-bound is being used, this can lead to extra nodes being checked, but it simplifies implementation and provides better correctness guarantees.

Inverted index

The LCA index in `sourmash` is an inverted index that stores a mapping from hashes in a collection of signatures to a list of IDs for signatures containing the hash. Despite the name, the list of signature IDs is not collapsed to the lowest common ancestor (as in `kraken`), and is calculated as needed by downstream methods using taxonomy information stored separately in the LCA index.

The mapping from hashes to signature IDs in the LCA index is an implicit representation of the original signatures used to build the index, and so returning the signatures is implemented by rebuilding the original signatures on-the-fly. Search in an LCA index matches the k -mers in the query to the list of signatures IDs containing them, using a counter data structure to sort results by number of hashes per signature ID. The rebuilt signatures are then returned as matches based on the signature ID, with containment or similarity to the query calculated against the rebuilt signatures.

`mash screen` [ondov mash 2019?] has a similar index, but it is constructed on-the-fly using the distinct hashes in a sketch collection as keys, and values are counters initially set to zero. As the query is processed, matching hashes have their counts incremented, and after all hashes in the query are processed then all the sketches in the collection are checked in the counters to quantify the containment/similarity of each sketch in the query. The LCA index uses the opposite approach, opting to reconstruct the sketches on-the-fly.

References

1. **On the resemblance and containment of documents**
AZ Broder
Institute of Electrical and Electronics Engineers (IEEE) (2002-11-23) <https://doi.org/fqk7hr>
DOI: [10.1109/sequen.1997.666900](https://doi.org/10.1109/sequen.1997.666900)
2. **IMPROVING MIN HASH VIA THE CONTAINMENT INDEX WITH APPLICATIONS TO METAGENOMIC ANALYSIS**
David Koslicki, Hooman Zabeti
Cold Spring Harbor Laboratory (2017-09-04) <https://doi.org/ghvn6z>
DOI: [10.1101/184150](https://doi.org/10.1101/184150)
3. **Mash Screen: high-throughput sequence containment estimation for genome discovery**
Brian D Ondov, Gabriel J Starrett, Anna Sappington, Aleksandra Kostic, Sergey Koren, Christopher B Buck, Adam M Phillippy
Genome Biology (2019-11-05) <https://doi.org/ghtqmb>
DOI: [10.1186/s13059-019-1841-x](https://doi.org/10.1186/s13059-019-1841-x) · PMID: [31690338](https://pubmed.ncbi.nlm.nih.gov/31690338/) · PMCID: [PMC6833257](https://pubmed.ncbi.nlm.nih.gov/PMC6833257/)
4. **Comparative metagenomic and rRNA microbial diversity characterization using archaeal and bacterial synthetic communities**
Migun Shakya, Christopher Quince, James H Campbell, Zamin K Yang, Christopher W Schadt, Mircea Podar
Environmental Microbiology (2013-06) <https://doi.org/f42ccr>
DOI: [10.1111/1462-2920.12086](https://doi.org/10.1111/1462-2920.12086) · PMID: [23387867](https://pubmed.ncbi.nlm.nih.gov/23387867/) · PMCID: [PMC3665634](https://pubmed.ncbi.nlm.nih.gov/PMC3665634/)
5. **Omega: an Overlap-graph de novo Assembler for Metagenomics**
Bahlul Haider, Tae-Hyuk Ahn, Brian Bushnell, Juanjuan Chai, Alex Copeland, Chongle Pan
Bioinformatics (2014-10) <https://doi.org/f6kt42>
DOI: [10.1093/bioinformatics/btu395](https://doi.org/10.1093/bioinformatics/btu395) · PMID: [24947750](https://pubmed.ncbi.nlm.nih.gov/24947750/)
6. **metaSPAdes: a new versatile metagenomic assembler**
Sergey Nurk, Dmitry Meleshko, Anton Korobeynikov, Pavel A Pevzner
Genome Research (2017-05) <https://doi.org/f97jkv>
DOI: [10.1101/gr.213959.116](https://doi.org/10.1101/gr.213959.116) · PMID: [28298430](https://pubmed.ncbi.nlm.nih.gov/28298430/) · PMCID: [PMC5411777](https://pubmed.ncbi.nlm.nih.gov/PMC5411777/)
7. **Evaluating Metagenome Assembly on a Simple Defined Community with Many Strain Variants**
Sherine Awad, Luiz Irber, CTitus Brown
Cold Spring Harbor Laboratory (2017-07-03) <https://doi.org/ghvn6x>
DOI: [10.1101/155358](https://doi.org/10.1101/155358)
8. **Mash: fast genome and metagenome distance estimation using MinHash**
Brian D Ondov, Todd J Treangen, Páll Melsted, Adam B Mallonee, Nicholas H Bergman, Sergey Koren, Adam M Phillippy
Genome Biology (2016-06-20) <https://doi.org/gfx74q>
DOI: [10.1186/s13059-016-0997-x](https://doi.org/10.1186/s13059-016-0997-x) · PMID: [27323842](https://pubmed.ncbi.nlm.nih.gov/27323842/) · PMCID: [PMC4915045](https://pubmed.ncbi.nlm.nih.gov/PMC4915045/)
9. **Critical Assessment of Metagenome Interpretation—a benchmark of metagenomics software**
Alexander Sczyrba, Peter Hofmann, Peter Belmann, David Koslicki, Stefan Janssen, Johannes Dröge, Ivan Gregor, Stephan Majda, Jessika Fiedler, Eik Dahms, ... Alice C McHardy
Nature Methods (2017-10-02) <https://doi.org/gbzspt>
DOI: [10.1038/nmeth.4458](https://doi.org/10.1038/nmeth.4458) · PMID: [28967888](https://pubmed.ncbi.nlm.nih.gov/28967888/) · PMCID: [PMC5903868](https://pubmed.ncbi.nlm.nih.gov/PMC5903868/)

10. **Tutorial: assessing metagenomics software with the CAMI benchmarking toolkit**
Fernando Meyer, Till-Robin Lesker, David Koslicki, Adrian Fritz, Alexey Gurevich, Aaron E Darling, Alexander Sczyrba, Andreas Bremges, Alice C McHardy
Nature Protocols (2021-03-01) <https://doi.org/gh77rh>
DOI: [10.1038/s41596-020-00480-3](https://doi.org/10.1038/s41596-020-00480-3) · PMID: [33649565](https://pubmed.ncbi.nlm.nih.gov/33649565/)
11. **GitHub - richarddurbin/modimizer: a toolset for fast DNA read set matching and assembly using a new type of reduced kmer**
GitHub
<https://github.com/richarddurbin/modimizer>
12. **Finch: a tool adding dynamic abundance filtering to genomic MinHashing**
Roderick Bovee, Nick Greenfield
The Journal of Open Source Software (2018-02-01) <https://doi.org/gm85dx>
DOI: [10.21105/joss.00505](https://doi.org/10.21105/joss.00505)
13. **sourmash: a library for MinHash sketching of DNA**
C Titus Brown, Luiz Irber
The Journal of Open Source Software (2016-09-14) <https://doi.org/ghdrk5>
DOI: [10.21105/joss.00027](https://doi.org/10.21105/joss.00027)

Scaled MinHash sketches support efficient indexing for large-scale containment queries

CTB: Additional points to raise:

- in-memory representation of sketches may be too big (!!), goal here is on disk storage/low minimum memory for “extremely large data” situation.
- Also/in addition, want ability to do incremental loading of things.
- Note we are not talking here about situations where the indices themselves are too big to download.
- I think rename LCA to revindex. Or make up a new name.

We provide two index data structures for rapid estimation of containment in large databases. The first, the MinHash Bloom Tree (MHBT), is a specialization of the Sequence Bloom Tree [solomon fast 2016?], and implements a k -mer aggregative method with explicit representation of datasets based on hierarchical indices. The second is LCA, an inverted index into sketches, a color-aggregative method with implicit representation of the sketches.

We evaluated the MHBT and LCA databases by constructing and searching a GenBank snapshot from July 18, 2020, containing 725,331 assembled genomes (5,282 Archaea, 673,414 Bacteria, 6,601 Fungi 933 Protozoa and 39,101 Viral). MHBT indices were built with *scaled* = 1000, and LCA indices used *scaled* = 10000. Table 2 shows the indexing results for the LCA index, and Table 3 for the MHBT index.

Table 2: Results for LCA indexing, with *scaled* = 10000 and k = 21.

Domain	Runtime (s)	Memory (MB)	Size (MB)
Viral	57	33	2
Archaea	58	30	5
Protozoa	231	3	17

Domain	Runtime (s)	Memory (MB)	Size (MB)
Fungi	999	3	65
Bacteria	12,717	857	446

Table 3: Results for MHBT indexing, with $scaled = 1000$, $k = 21$ and internal nodes (Bloom Filters) using 10000 slots for storage.

Domain	Runtime (s)	Memory (MB)	Size (MB)
Viral	126	326	77
Archaea	111	217	100
Protozoa	206	753	302
Fungi	1,161	3,364	1,585
Bacteria	32,576	47,445	24,639

Index sizes are more affected by the number of genomes inserted than the individual *Scaled MinHash* sizes. Despite Protozoan and Fungal *Scaled MinHash* sketches being larger individually, the Bacterial indices are an order of magnitude larger for both indices since they contain two orders of magnitude more genomes.

Comparing between LCA and MHBT index sizes must account for their different scaled parameters, but as shown in Chapter 1 a *Scaled MinHash* with $scaled = 1000$ when downsampled to $scaled = 10000$ is expected to be ten times smaller. Even so, MHBT indices are more than ten times larger than their LCA counterparts, since they store extra caching information (the internal nodes) to avoid loading all the data to memory during search. LCA indices also contain extra data (the list of datasets containing a hash), but this is lower than the storage requirements for the MHBT internal nodes.

We next executed similarity searches on each database using appropriate queries for each domain. All queries were selected from the relevant domain and queried against both MHBT ($scaled = 1000$) and LCA ($scaled = 10000$), for $k = 21$.

Table 4: Running time in seconds for similarity search using LCA ($scaled = 10000$) and MHBT ($scaled = 1000$) indices.

	Viral	Archaea	Protozoa	Fungi	Bacteria
LCA	1.06	1.42	5.40	26.92	231.26
SBT	1.32	3.77	43.51	244.77	3,185.88

Table 5: Memory consumption in megabytes for similarity search using LCA ($scaled = 10000$) and MHBT ($scaled = 1000$) indices.

	Viral	Archaea	Protozoa	Fungi	Bacteria
LCA	223	240	798	3,274	20,926
SBT	163	125	332	1,656	2,290

Table 4 shows running time for both indices. For small indices (Viral and Archaea) the LCA running time is dominated by loading the index in memory, but for larger indices the cost is amortized due to the faster running times. This situation is clearer for the Bacteria indices, where the LCA search completes in 3 minutes and 51 seconds, while the SBT search takes 54 minutes.

When comparing memory consumption, the situation is reversed. Table [5](#) shows how the LCA index consistently uses twice the memory for all domains, but for larger indices like Bacteria it uses as much as 10 times the memory as the MHBT index for the same data.

For both runtime and memory consumption, it is worth pointing that the LCA index is a tenth of the data indexed by the MHBT. This highlights the trade-off between speed and memory consumption for both approaches, especially for larger indices.

Notes: * new genomes can be added quickly to SBT.