

# Strain Level Microbial Detection and Quantification with Applications to Single Cell Metagenomics

Kaiyuan Zhu<sup>1,2</sup>, Welles Robinson<sup>2,3</sup>, Alejandro A. Schäffer<sup>2</sup>, Junyan Xu<sup>2</sup>, Eytan Ruppin<sup>2</sup>,  
A. Funda Ergun<sup>1</sup>, Yuzhen Ye<sup>1</sup>, and S. Cenk Sahinalp<sup>2</sup>

<sup>1</sup>Department of Computer Science, Indiana University, Bloomington, IN, USA

<sup>2</sup>Cancer Data Science Laboratory, National Cancer Institute, National Institutes of Health,  
Bethesda, MD, USA

<sup>3</sup>Center for Bioinformatics and Computational Biology, University of Maryland, College  
Park, MD, USA

## Abstract

The identification and quantification of microbial abundance at the species or strain level from sequencing data is crucial for our understanding of human health and disease. Existing approaches for microbial abundance estimation either use accurate but computationally expensive alignment-based approaches for species-level estimation or less accurate but computationally fast alignment-free approaches that fail to classify many reads accurately at the species or strain-level.

Here we introduce CAMMiQ, a novel combinatorial solution to the microbial identification and abundance estimation problem, which performs better than the best used tools on simulated and real datasets with respect to the number of correctly classified reads (i.e., specificity) by an order of magnitude and resolves possible mixtures of similar genomes.

As we demonstrate, CAMMiQ can better distinguish between single cells deliberately infected with distinct *Salmonella* strains and sequenced using scRNA-seq reads than alternative approaches. We also demonstrate that CAMMiQ is also more accurate than the best used approaches on a variety of synthetic genomic read data involving some of the most challenging bacterial genomes derived from NCBI RefSeq database; it can distinguish not only distinct species but also closely related strains of bacteria.

The key methodological innovation of CAMMiQ is its use of arbitrary length, doubly-unique substrings, i.e. substrings that appear in (exactly) two genomes in the input database, instead of fixed-length, unique substrings. To resolve the ambiguity in the genomic origin of doubly-unique substrings, CAMMiQ employs a combinatorial optimization formulation, which can be solved surprisingly quickly. CAMMiQ's index consists of a sparsified subset of the shortest unique and doubly-unique substrings of each genome in the database, within a user specified length range and as such it is fairly compact. In short, CAMMiQ offers more accurate genomic identification and abundance estimation than the best used alternatives while using similar computational resources.

**Availability:** <https://github.com/algocancer/CAMMiQ>

## 1 Introduction

Advances in high throughput sequencing (HTS) have made it possible to generate millions of short reads in a few hours. An increased appreciation for the importance of microbes in human health and disease has prompted the generation of many metagenomic HTS datasets. For example, whole genome shotgun sequencing of the gut microbiome by the Human Microbiome Project provided important insight into the function and diversity of the human gut microbiome [1]. Similarly, the increase in HTS of human tissues represents another enormous source of metagenomic data because many of these human HTS datasets include reads from tissue-resident microbes, which have been shown to play an important role in many aspects of human disease, including tumorigenesis and the tumor response to therapy [2, 3, 4, 5, 6, 7].

Early approaches for analyzing metagenomic sequencing data were alignment based; reads were primarily searched in GenBank [8] through BLASTN [9] or custom built aligners such as GATK PathSeq (which, in

return, is based on popular alignment tools) [10]. Unfortunately, the recent growth of HTS data and reference databases has made read search and alignment using BLASTN or PathSeq computationally infeasible. For example, a recent study which reported that microbial reads from tumors sequenced by The Cancer Genome Atlas (TCGA) can be used to build a classifier for cancer type [11] needed an “alignment-free” approach (Kraken [12]) due to the large number of samples analyzed. Kraken and other recent computational methods aim to identify and quantify specific genus in a metagenomic sample faster through a number of algorithmic techniques are summarized in [13]. Even though these methods address the problem of computational scalability, available alignment-free tools do not match the accuracy of alignment-based tools. For example, another recent paper on microbial reads from single cell RNA-seq (scRNA-seq) datasets to distinguish cell type specific intracellular microbes from extracellular and contaminating microbes [14] had to use PathSeq. This was because of the relatively small number of microbial reads per cell, to classify as many reads as possible at the genus and species levels - the resolution of available alignment-free methods were too limited for this task. The distinct approaches taken by these two studies illustrate the trade-offs inherent in the alignment-based and alignment-free methods to metagenomic identification and quantification.

One additional way to achieve speedup is through the reduction of the size of the sequence database used. Methods in this direction may align reads only to *marker genes*, a relatively small collection of clade-specific, single-copy genes, instead of the full reference genomes [15, 16], through the use of available read mapping techniques. Such methods need to obtain the collection of marker genes by employing information beyond what is offered by the database itself - which is not always possible. Furthermore, since these methods identify and use only a handful of marker genes on each genome, many of the reads in the HTS data can not be utilized, implying low specificity. As a consequence, species with low abundance within the sample may be difficult to identify (let alone correctly quantify) because variation in HTS coverage may result in a few or no reads originating from the marker genes.

Alignment-free methods have been available in the context of string comparison for a long time [17, 18, 19] and they have found applications in bioinformatics workflows before the emergence of HTS [20, 21]. Applications of alignment-free methods to metagenomic HTS data typically rely on  $k$ -mer “matches” to return a taxonomic assignment for every read. These applications either assign a read to the lowest taxonomic rank possible (determined by the specificity of the read’s  $k$ -mers) [22, 12, 23, 24], or to a pre-determined genus, species, or strain taxonomic level [25, 26]. In contrast to marker gene based methods,  $k$ -mer based applications can usually assign each given read to the correct genus. As the value of  $k$  becomes larger, more reads can be assigned a unique label at the desired taxonomic rank, however, with growing  $k$ , the space requirement of these methods grows very quickly - which, in the worst case, can imply a factor  $k$  increase on the space requirements of the original database. The large memory footprint to maintain the entire  $k$ -mer profile of each species can be reduced through hashing or subsampling the  $k$ -mers [27, 28]; however this would result in loss of accuracy. In addition to methods based on exact  $k$ -mer counts, it is also possible to assign metagenomic reads to bacterial genomes by employing species-specific sequence features (e.g. short  $k$ -mer distribution or GC content) [29, 30, 31, 32, 33], although methods that employ this approach are typically not very accurate at species level or (especially) at strain level assignment.

The ability to identify and quantify distinct pathogenic strains in a microbial sample has many applications in diagnosing and treating infectious diseases [34]. In particular, mixed infections that are caused by multiple strains of a bacterial species [35, 36], have been described for tens of bacterial species [37]. It is estimated that a high fraction of *M. tuberculosis* and *Staphylococcus aureus* patients are infected with multiple strains [34], each with distinct different drug susceptibility and antibiotic resistance profiles. Clearly accurate analysis of microbial sequence data at the strain level would be highly helpful for identifying mixed infections. Additionally, strain-level analysis can be important for distinguishing pathogenic strains from non-pathogenic strains [38] and for tracking food-borne pathogens [39].

As have been summarized above, neither the alignment free,  $k$ -mer based methods, nor the marker-gene based tools take into consideration the distribution of the reads over the genome of the species they are assigned. In fact, the alignment-free methods do not even take into account the length of the genomes, but are based on (unnormalized) read counts. In reality, provided that the sequence data to be analyzed is genomic, the distribution of reads generated by high throughput sequencing from a given species or strain should be roughly uniform. This principle is fundamental to a number of *isoform abundance estimation* methods, which aim to solve a very similar problem [40, 41, 42]. Interestingly, this constraint is underutilized in the context of metagenomic analysis. One exception is the network flow based approach utilized

by, e.g. [43], which (implicitly) establishes a reference guided assembly of the reads into the genomes of the species involved. As such it is quite accurate but is very slow. Another method in this direction considers the uniformity of coverage across  $k$ -mers with each genome to reduce false positive calls [44], which improves the running time at a moderate loss of accuracy.

In contrast to the metagenomic species identification and quantification methods summarized above, there are also tools to determine the likely presence of a long genomic sequence (e.g. the complete or partial genome of a bacterial species) in a given metagenomic sample [45, 46, 47, 48]. Even though these tools solve an entirely different problem, methodologically they are similar to the  $k$ -mer based species identification and quantification tools such as [12, 26] in the sense that they build a succinct index on the database (this time comprised of the the metagenomic read collection) and query this index without explicit alignment. And because of their design parameters, these tools can not perform abundance estimation for a given species.

**Our Contributions.** In this paper we describe CAMMiQ (Combinatorial Algorithms for Metagenomic Microbial Quantification), a new computational method to maintain/manage a collection of  $m$  (bacterial) genomes  $\mathcal{S} = \{s_1, \dots, s_m\}$ , each assembled into one or more strings/contigs, representing a species, a particular strain of a species, or any other taxonomic rank. CAMMiQ’s data structure can answer queries of the following form: given a query set  $\mathcal{Q}$  of HTS reads obtained from a mixture of genomes or transcriptomes, each from  $\mathcal{S}$ , identify the genomes in  $\mathcal{Q}$ , and, in case the reads are genomic, compute their relative abundances. Our data structure is not only very efficient in terms of its querying time, but is also shown to be very accurate, through simulations. The key novel feature of our data structure is its utilization of substrings that are present in at most  $c$  genomes ( $c > 1$ ) in  $\mathcal{S}$ . There are alignment free methods that utilize unique  $k$ -mers in genomes for metagenomic analysis [12, 26, 44] already. However, our data structure is the first to consider those substrings that are present in two or possibly more genomes, for increasing the proportion of reads it can utilize and thus improving sensitivity. Because of this novel feature, our data structure can accurately identify genomes at subspecies/strain level.

In order to assign each read in  $\mathcal{Q}$  that includes an “almost-unique” substring (i.e. present in at most  $c$  genomes) to a genome, our data structure solves an integer linear program (ILP) - that simultaneously infers which genomes are present in  $\mathcal{Q}$  and, (if the reads are genomic) their relative abundances. Specifically, the objective of the ILP is to identify a set of genomes, in each of which the coverage of the almost-unique genomic substrings is (approximately) uniform.

One novel feature of our data structure is its use of *shortest* substrings (present in at most  $c$  genomes) - rather than fixed length “ $k$ -mers”.<sup>1</sup> This feature also increases the number of reads utilized by our data structure since some reads may include no almost-unique  $k$ -mer but instead longer substrings that are almost-unique. It also reduces the number of substrings and thus features (not to be confused with the number of reads) to be handled by the ILP: this is because, rather than considering multiple overlapping (almost-)unique  $k$ -mers, our data structure uses a single shorter substring shared by all of them. The ILP formulation further reduces the number of almost-unique substrings it considers by maintaining a maximally sparsified set of substrings that guarantee any potential read from a genome in  $\mathcal{S}$  that is almost-unique includes one such substring. In the remainder of the paper, we set the value of  $c$  to 2, so as to consider not only unique but also “doubly-unique” substrings of the genomes in  $\mathcal{S}$ . Even though this choice turned out to be sufficiently powerful for the datasets we experimented with, it is easy to generalize our approach to  $c > 2$  (i.e. to triply-unique, etc. substrings).

As a final contribution, we provide sufficient conditions to identify and quantify genomes in a query correctly, through the use of unique substrings/ $k$ -mers, provided the reads are error-free. Although this is a purely theoretical result, to the best of our knowledge it is the first of its kind for metagenomic data analysis, and is valid for CAMMiQ for the case  $c = 1$  and other unique substring based methods such as CLARK and KrakenUniq. CAMMiQ’s use with  $c = 2$  is primarily advised for cases where these conditions are not satisfied.

We show that CAMMiQ is not only faster but also more accurate than GATK PathSeq on scRNA-seq data obtained from monocyte-derived dendritic cells (moDCs) infected with distinct *Salmonella* strains analyzed in [14]. Then we demonstrate the comparative advantage of CAMMiQ against some of the best used

---

<sup>1</sup>There are also the so-called “gapped”  $k$ -mers used in a variety of applications such as functional classification of metagenomic reads [49]; however not only these applications are different from that of CAMMiQ, but also, they use fixed length substrings. Finally, there are alternative approaches using, e.g. linked-reads [50]; CAMMiQ focuses on standard Illumina reads to be able to analyze a wealth of readily available sequencing data sets.

metagenomic tools on synthetic genomic HTS data from “challenging” microbial strains we derived from the NCBI RefSeq database.

## 2 Algorithmic Formulation

The input to CAMMIQ is a set of  $m$  genomes  $\mathcal{S} = \{s_i\}_{i=1}^m$ , not necessarily all from the same taxonomic level (each genome here may be associated with a genus, species, subspecies or strain) to be indexed. Although we describe CAMMIQ for the case where each  $s_i \in \mathcal{S}$  is a single string, we do not assume that the genomes are fully assembled into a single contig; rather the string representing a genome could simply be a concatenation of all contigs from species  $i$  and their reverse complements with a special symbol  $\$i$  between each contig. We call  $\mathcal{S}$  the input *database* and  $i \in \{1, \dots, m\}$  the *genome ID* of string  $s_i$ .

A query for CAMMIQ involves a set of reads  $\mathcal{Q} = \{r_j\}_{j=1}^n$  representing a metagenomic mixture. For simplicity we describe CAMMIQ for reads of length  $L$ , however our data structure can handle reads of varying length. Given  $\mathcal{Q}$ , the goal of CAMMIQ is to identify a set of genomes  $\mathcal{A} = \{s_1, \dots, s_a\} \subset \mathcal{S}$  and their respective abundances  $p_1, \dots, p_a$  which “best explain”  $\mathcal{Q}$ . This is achieved by assigning (select) reads  $r_j$  to genomes  $s_i$  such that the implied coverage of each genome  $s_i \in \mathcal{A}$  is uniform across  $s_i$ , with  $p_i$  as the mean.

In its simplest form, CAMMIQ builds a succinct index for input database  $\mathcal{S}$ , so as to handle queries in the following form. Given  $\mathcal{Q} = \{r_1, \dots, r_n\}$  find the smallest set of genomes  $\mathcal{A}_1 = \{s_1, \dots, s_{a_1}\} \subseteq \mathcal{S}$ , such that the set of *unique*  $L$ -mers in  $\mathcal{A}_1$  includes all reads in  $\mathcal{Q}$  that are identical to a unique  $L$ -mer in  $\mathcal{S}$ . For any  $s_i \in \mathcal{S}$ , call an  $L$ -mer in  $s_i$  unique in  $\mathcal{S}$  if it does not appear in any other  $s_{i'}$ . We denote by  $U_{i,L}$ , the set of all unique  $L$ -mers in  $s_i$ . Then  $\mathcal{A}_1 = \{s_i \in \mathcal{S} \mid \mathcal{Q} \cap U_{i,L} \neq \emptyset\}$ .

The above query type may not be powerful enough to identify in  $\mathcal{Q}$  those genomes in  $\mathcal{S}$  that are very similar (with respect to sequence composition) to other genomes and thus do not include many unique  $L$ -mers. Such a genome  $s_i$  would be especially problematic if it is low in abundance, since the chances of  $\mathcal{Q}$  to include any unique  $L$ -mer from  $U_{i,L}$  will be low. As a solution to this problem, CAMMIQ also features a second type of query, which is more general since it involves both unique and *doubly-unique*  $L$ -mers of genomes  $s_i$ . We call an  $L$ -mer doubly-unique in  $\mathcal{S}$  if it appears in exactly two genomes  $s_i$  and  $s_{i'}$  and denote the set of doubly-unique  $L$ -mers in  $s_i$  by  $D_{i,L}$ . Given query set  $\mathcal{Q}$ , CAMMIQ’s more general query asks to compute  $\mathcal{A}_2$ , the smallest subset of  $\mathcal{S}$  which include all reads in  $\mathcal{Q}$  identical to unique and doubly-unique  $L$ -mers in  $\mathcal{S}$ . Note that, necessarily  $\mathcal{A}_1 \subseteq \mathcal{A}_2$ . Also note that, by involving doubly-unique  $L$ -mers, this more general query may have a chance to capture those genomes in  $\mathcal{Q}$  which have low abundance and have a few unique  $L$ -mers.

Unfortunately, the index structure necessary to support the second type of queries is larger than that for the first type of queries. Furthermore, this query type could still produce inaccurate results in the presence of read errors. For handling read errors (at noise rates commonly observed in Illumina data) CAMMIQ finally features a third type of query, which is even more general since it involves the “shortest” unique and doubly-unique substrings of reads in  $\mathcal{Q}$ . We call a substring of a genome  $s_i \in \mathcal{S}$  a shortest unique substring, if it is unique to  $s_i$ , has length in the range  $[L_{\min}, L_{\max}]$ , and has no substring that is a shortest unique substring; this definition can be extended to shortest doubly-unique substrings as well. This query thus asks to compute the smallest subset  $\mathcal{A}_3$  of  $\mathcal{S}$  which include these substrings, with the constraint that the “coverage” of these substrings in each genome  $s_i \in \mathcal{A}_3$  is “uniform”. The query also asks to compute the relative abundance of each genome in  $\mathcal{A}_3$  as will be described later.

By considering shorter unique substrings, our most general query has a higher chance of observing them within  $\mathcal{Q}$  (since a substring of length  $L' < L$  is included in  $L - L' + 1$   $L$ -mers in the genome). This query may also imply that a smaller fraction of unique or doubly-unique substrings in  $\mathcal{Q}$  would be subject to read errors (since it is likely that each read may include more than one such unique substring and it is likely that at least one such unique substring will be error free). Furthermore the index structure necessary to maintain the entire set of shortest unique and doubly-unique substrings of genomes in  $\mathcal{S}$  is fairly large. As a result of this, we build the index on a maximally sparsified set of shortest unique and doubly-unique substrings of each  $s_i \in \mathcal{S}$  which ensure that each read in  $\mathcal{Q}$  that can be attributed to at most two genomes includes at least one such substring (see the Supplementary Methods for more details).

An index on this sparsified set of shortest unique and doubly-unique substrings is sufficiently powerful for CAMMIQ to answer all three types of queries, i.e. it can efficiently compute the sets  $\mathcal{A}_1$ ,  $\mathcal{A}_2$  and  $\mathcal{A}_3$ . For all

three query types CAMMiQ first identifies for each read  $r_j$  all unique and doubly-unique substrings it includes; it then assigns  $r_j$  to the one or two genomes from which these substrings can originate.<sup>2</sup> To compute  $\mathcal{A}_1$  CAMMiQ can simply return the collection of genomes which receive at least one read assignment. To compute  $\mathcal{A}_2$  CAMMiQ needs to solve the “set cover” problem, or more precisely, its dual, the “hitting set” problem where genomes form sets and indexed strings that appear in query reads form the items to be covered. To compute  $\mathcal{A}_3$  CAMMiQ solves the combinatorial optimization problem that asks to minimize the variance among the number of reads assigned to each indexed substring of each genome - the solution indicates the set of genomes in  $\mathcal{A}_3$  along with their respective abundances.

Details on the composition as well as the construction process for CAMMiQ’s index are discussed in Section 2.1. The query processing of CAMMiQ, which proceeds in two stages are discussed in Sections 2.2 and 2.3: The first stage assigns reads to specific genomes - which is sufficient for computing sets  $\mathcal{A}_1$  and  $\mathcal{A}_2$ . See Section 2.2 for the criteria we use for assigning a read to a genome, based on the indexed substrings it includes. The second stage introduces the combinatorial optimization formulation to compute  $\mathcal{A}_3$  as a response to the most general query type. See Section Section 2.3 for details.

## 2.1 Index Construction

In order to respond to the three types of queries described above, we preprocess unique and doubly-unique substrings of genomes in  $\mathcal{S}$  to form an index structure as follows. Let  $\mathcal{U} = \cup_{i=1}^m \mathcal{U}_i$  and  $\mathcal{D} = \cup_{i=1}^m \mathcal{D}_i$  where  $\mathcal{U}_i$  and  $\mathcal{D}_i$  are the substrings from genome  $s_i$  whose lengths are within the range than  $[L_{\min}, L_{\max} \leq L]$ , such that each  $u \in \mathcal{U}$  is present in at most one genome and each  $d \in \mathcal{D}$  is present in at most two genomes in  $\mathcal{S}$ . See below for a detailed definition for the uniqueness of a substring. If our goal is to compute  $\mathcal{A}_1$  or  $\mathcal{A}_2$  only (i.e. respond to the first and second type of queries), we can simply set  $L_{\min} = 0, L_{\max} = L$  and maintain the corresponding substring collections  $\mathcal{U}$  and  $\mathcal{D}$  in our index. For computing  $\mathcal{A}_3$  (i.e. responding to the third type of queries - in addition to the first and second types), we actually use the length constrained definitions of substrings for computing  $\mathcal{U}$  and  $\mathcal{D}$ . As mentioned earlier, we then sparsify  $\mathcal{U}$  and  $\mathcal{D}$  as much as possible by maintaining only one representative substring among those that are in close proximity within a genome, and discarding the rest (see “Subsampling unique substrings” below for details). We start with formal definitions and some notation.

**Notation and definitions.** Let  $s = s_1\$1 \circ \dots \circ s_m\$m$  denote the string obtained by concatenating the input reference genomes  $s_i \in \mathcal{S}$  and let  $M = |s| = \sum_i |s_i|$  denote its length. A substring of  $s$  is a string in the form  $s[l : r] = s[l]s[l+1] \dots s[r]$ . With a slight abuse of notation we denote by  $s_i[l : r]$  not the actual substring of  $s_i$  including its  $l^{th}$  to  $r^{th}$  symbols, but rather the substring of  $s$  including its  $l^{th}$  to  $r^{th}$  symbols, with the provision that all these symbols are within the representation of  $s_i$  in  $s$ . We denote by an  $\ell$ -mer a string of length  $\ell$ . The suffix of  $s$  that starts at position  $i$  is denoted  $\text{su}f[i] = s[i, \dots, M]$ . In what follows, we use the *generalized enhanced suffix array* of  $s$  which is composed of three parts. (i) The *suffix array* SA of  $s$ , which is comprised of the positions  $1, 2, \dots, M$ , sorted in increasing lexicographical order of the corresponding suffixes  $\text{su}f[i], i = 1, 2, \dots, M$ . That is,  $\text{SA}[i] = j$  indicates that  $\text{su}f[j]$  is the  $i$ -th smallest suffix in lexicographical order. In addition, we denote  $\text{SA}^{-1}[j] = i$  if  $\text{SA}[i] = j$ . (ii) The *longest common prefix array*, LCP, contains in its  $i$ -th position the length of the longest common prefix of  $\text{su}f[\text{SA}[i]]$  and  $\text{su}f[\text{SA}[i-1]]$ , for  $2 \leq i \leq M$  (and  $\text{LCP}[1] = 0$ ). (iii) Finally, the *generalized suffix array* GSA contains the genome ID of each suffix  $\text{su}f[\text{SA}[i]]$ . All of the above arrays can be constructed in linear time: the first data structure that can be constructed in linear time, with the ability to determine whether a given substring is unique to a “document” (i.e. a genome in our context) in a collection of documents, and compute the shortest unique substring of a document that ends in a particular position, in time proportional to the substring length is the augmented suffix tree of Matias et al. [51]. Once an augmented suffix tree is computed, it can be trivially reduced to the above described enhanced suffix array in  $O(M)$  time - which can also be constructed without the use of suffix trees to achieve a constant factor improvement in memory [52, 53].

We denote by  $u_i(l, r)$  the substring  $s_i[l : r]$  that is *unique* to genome  $s_i$ ; formally, this indicates that there exists no substring  $u_j(l', r')$  on any genome  $s_j \neq s_i$  such that  $u_i(l, r) = u_j(l', r')$ . We call  $u_i(l, r)$  a shortest unconstrained unique substring, if none of its substrings are unique. Similarly, we denote by  $d_i(l, r)$  the

<sup>2</sup>We explain how CAMMiQ resolves ambiguity on reads that include substrings which suggest conflicting assignments later in the paper.



substring  $s_i[l : r]$  that is *doubly-unique* to genome  $s_i$  and one other genome, say  $s_j$ ; formally, this indicates that there is exactly one genome  $s_j$  which includes the substring  $d_j(l', r')$ , i.e.,  $s_i[l : r] = s_j[l' : r']$  for some  $l', r'$ . Clearly, any superstring of a unique substring is still unique and any superstring of a doubly-unique substring is either unique or doubly-unique. We call  $d_i(l, r)$  a shortest unconstrained doubly-unique substring of  $s_i$  and some other genome  $s_j$ , if none of its substrings are doubly-unique.

For our purposes, we need to constrain the shortest unique and doubly-unique substrings with length upper bound  $L_{\max}$  and lower bound  $L_{\min}$ . Under these constraints, we call any shortest unconstrained unique substring  $u_i(l, r)$  a shortest unique substring if  $L_{\max} \geq r - l + 1 > L_{\min}$ . We also call a unique substring  $u_i(l, r)$  a shortest unique substring if  $r - l + 1 = L_{\min}$ . Similarly, we call any shortest unconstrained doubly-unique substring  $d_i(l, r)$  a shortest doubly-unique substring if  $L_{\max} \geq r - l + 1 > L_{\min}$ . Again, we call a doubly-unique substring  $u_i(l, r)$  a shortest doubly-unique substring if  $r - l + 1 = L_{\min}$  as well. We say an  $L$ -mer  $s_i[l : l + L - 1]$  includes a unique substring  $s_i[l' : r']$ , or, conversely, a unique substring  $s_i[l' : r']$  covers an  $L$ -mer  $s_i[l : l + L - 1]$  if  $l' \geq l$  and  $r' \leq l + L - 1$ . As such, we call an  $L$ -mer unique if it includes a unique substring. We can generalize these definitions to the notion of an  $L$ -mer including a doubly-unique substring, or conversely, a doubly-unique substring covering an  $L$ -mer, and thus making the  $L$ -mer itself doubly-unique - provided that it is not unique.

**Algorithmic framework to compute shortest unique substrings.** It is quite simple to compute the shortest unique and doubly-unique substrings in  $\mathcal{S}$  in  $O(M)$  time by using the augmented suffix tree described in [51]. A similar running time can also be achieved through the use of a suffix array, as discussed by [54] for a single document (i.e. genome). We slightly generalize this to handle multiple genomes as follows. The key observation we use is that given a position  $l$ , the shortest unique or doubly-unique substring of  $s_i$  that starts at  $l$  (i.e.  $u_i(l, r)$  or  $d_i(l, r)$ ) is the shortest unique, or respectively doubly-unique prefix of  $\text{suf}[l]$ . In this way the problem can be reduced to searching for the longest common prefix of  $\text{suf}[l]$  with any other suffix from another genome (i.e., any genome with  $\text{ID} \neq \text{GSA}[\text{SA}^{-1}[l]]$ ) for each  $1 \leq l \leq M$ . Let  $\text{lcp}(x, y)$  denote the longest common prefix of two suffixes  $x$  and  $y$ ; then we define:

$$\text{LCP}_u[i] = \max_{1 \leq j \leq M; \text{GSA}[j] \neq \text{GSA}[i]} \text{lcp}(\text{suf}[\text{SA}[i]], \text{suf}[\text{SA}[j]]) \quad (1)$$

and

$$\text{LCP}_d[i] = \max_{1 \leq j \leq M; \text{GSA}[j] \neq \text{GSA}[i], \text{GSA}[j']} \text{lcp}(\text{suf}[\text{SA}[i]], \text{suf}[\text{SA}[j]]) \quad (2)$$

where  $j'$  indicates a suffix  $\text{suf}[j']$  with  $\text{GSA}[j'] \neq \text{GSA}[i]$ , which maximizes  $\text{lcp}(\text{suf}[\text{SA}[i]], \text{suf}[\text{SA}[j']])$ ; note that any value of  $j'$  that maximizes  $\text{lcp}(\text{suf}[\text{SA}[i]], \text{suf}[\text{SA}[j']])$  will imply the same value for  $\text{LCP}_d[i]$ . CAMMiQ maintains an array SU (of length  $M$ ) such that  $\text{SU}[r] = l$  if  $u_i(l, r)$  is a shortest unique substring. In order to compute SU, each of its entries is initially set to 0 and for each  $i = 1, \dots, M$ , one entry of SU is updated as

$$\text{SU}[\text{SA}[i] + \text{LCP}_u[i]] \leftarrow \max\{\text{SU}[\text{SA}[i] + \text{LCP}_u[i]], \text{SA}[i]\} \quad (3)$$

Similarly, CAMMiQ maintains an array SD (again of length  $M$ ) such that  $\text{SD}[r] = l$  if  $d_i(l, r)$  is a shortest doubly-unique substring. Again each entry of SD is initially set to 0 and then for each  $i = 1, \dots, M$ , one entry of SD is updated as

$$\text{SD}[\text{SA}[i] + \text{LCP}_d[i]] \leftarrow \max\{\text{SD}[\text{SA}[i] + \text{LCP}_d[i]], \text{SA}[i]\}. \quad (4)$$

See Supplementary Section 5.1 for further details.

**Computing  $\text{LCP}_u$  and  $\text{LCP}_d$ .** Given  $\text{GSA}[i_1, \dots, i_2]$ , a subarray of GSA, let  $d_{\text{GSA}}(i_1, i_2)$  be the number of distinct genomes the entries in this subarray belong to, i.e.  $d_{\text{GSA}}(i_1, i_2) = |\{\text{GSA}[i_1], \dots, \text{GSA}[i_2]\}|$ . We can now compute  $\text{LCP}_u$  and  $\text{LCP}_d$  in linear time as follows.

$$\text{LCP}_u[i] = \max \begin{cases} \min_{i^- < x \leq i} \text{LCP}[x], \text{ where } i^- = \max\{1 \leq i' < i\}, \text{ s.t. } d_{\text{GSA}}(i', i) \geq 2 \\ \min_{i < x \leq i^+} \text{LCP}[x], \text{ where } i^+ = \min\{i < i' \leq M\}, \text{ s.t. } d_{\text{GSA}}(i, i') \geq 2 \end{cases} \quad (5)$$

$$\text{LCP}_d[i] = \min \left\{ \begin{array}{l} \max \left\{ \begin{array}{l} \min_{i^- < x \leq i} \text{LCP}[x], \text{ where } i^- = \max\{1 \leq i' < i\}, \text{ s.t. } d_{\text{GSA}}(i', i) \geq 2 \\ \min_{i < x \leq i^{2+}} \text{LCP}[x], \text{ where } i^{2+} = \min\{i < i' \leq M\}, \text{ s.t. } d_{\text{GSA}}(i, i') \geq 3 \end{array} \right. \\ \max \left\{ \begin{array}{l} \min_{i^{2-} < x \leq i} \text{LCP}[x], \text{ where } i^{2-} = \max\{1 \leq i' < i\}, \text{ s.t. } d_{\text{GSA}}(i', i) \geq 3 \\ \min_{i < x \leq i^+} \text{LCP}[x], \text{ where } i^+ = \min\{i < i' \leq M\}, \text{ s.t. } d_{\text{GSA}}(i, i') \geq 2 \end{array} \right. \end{array} \right. \quad (6)$$

Note that to introduce a minimum length constraint  $L_{\min}$  on unique and doubly-unique substrings, each  $\text{LCP}_u[i]$  is (re)set to  $\max\{L_{\min} - 1, \text{LCP}_u[i]\}$  and respectively each  $\text{LCP}_d[i]$  is (re)set to  $\max\{L_{\min} - 1, \text{LCP}_d[i]\}$ . Then, to ensure that each shortest doubly-unique substring occurs in exactly two genomes (and not one), we set  $\text{LCP}_d[i] = \infty$  in case the above procedure ends up with  $\text{LCP}_d[i] = \text{LCP}_u[i]$ . See Supplementary Section 5.2 for the proof of correctness and a running time analysis for the computation of  $\text{LCP}_u$  and  $\text{LCP}_d$ .

**Subsampling unique substrings.** As defined above, the array  $\text{SU}$  maintains the indices of all shortest unique substrings; some of these substrings may have large overlaps with others and thus are redundant in assessing the uniqueness of a read in the query. Let  $\mathcal{U}_i$  be the collection of all unique substrings on genome  $s_i$ . Then, in order to reduce the index size, **CAMMiQ** aims to compute a subset  $\mathcal{U}'_i$  of  $\mathcal{U}_i$ , consisting of the minimum number of shortest unique substrings that cover every unique  $L$ -mer on  $s_i$ . **CAMMiQ** also aims to compute a subset  $\mathcal{D}'_i$  of  $\mathcal{D}_i$ , consisting of the minimum number of shortest doubly-unique substrings that cover every doubly-unique  $L$ -mer on  $s_i$ . This is all done by greedily maintaining only the rightmost shortest unique or doubly-unique substring in any  $L$ -mer of a genome in  $\mathcal{S}$ . In the remainder of the paper we denote the number of unique substrings in subset  $\mathcal{U}'_i$  by  $nu_i$  ( $= |\mathcal{U}'_i|$ ) and respectively the number of doubly-unique substrings in subset  $\mathcal{D}'_i$  by  $nd_i$  ( $= |\mathcal{D}'_i|$ ); we denote the number of unique  $L$ -mers on  $s_i$  by  $nu_i^L$  and respectively the number of doubly-unique  $L$ -mers on  $s_i$  by  $nd_i^L$ . As we prove in Supplementary Section 5.3, this greedy strategy we employ can indeed obtain the minimum number of shortest unique substrings to cover each unique  $L$ -mer, provided that each substring in  $\mathcal{U}_i$  occurs only once  $s_i$ .

**Index Organization.** Consider the set of unique substrings Let  $h = \min_{u_i \in \mathcal{U}} |u_i|$  be the minimum length of all shortest unique substrings ( $h$  is automatically set to  $L_{\min}$  if the minimum length constraint is imposed). **CAMMiQ** maintains a hash table which maps a distinct  $h$ -mer  $w$  to a bucket containing all unique substrings  $u_i$  which have  $w$  as a prefix. Within each bucket, the remaining suffices of all unique substrings  $u_i$ , i.e.  $u_i[h+1 : |u_i|]$ , are maintained in a trie so that each leaf contains the corresponding genome ID. For each read  $r_j$  in the query, **CAMMiQ** considers each substring of length  $h$  and its reverse complement and computes its hash value (in time linear with  $L$  through Karp-Rabin fingerprinting [55]). If the substring has a match in the hash table, then **CAMMiQ** tries to extend the match until a matching unique substring is found, or finds no match. (Note that the processing for doubly-unique substrings is identical to that for unique substrings.) See Figure 1 for an overview of the index structure. Also see Section 2.2 below for the use of unique and doubly-unique substrings identified for each read to answer the query.

## 2.2 Query Processing Stage 1: Preprocessing the Reads

Given the index structure on the sparsified set of shortest unique and doubly-unique substrings of genomes in  $\mathcal{S}$ , we handle each query  $\mathcal{Q}$  in two stages. In Stage 1, we preprocess each read  $r_j \in \mathcal{Q}$  as follows. Consider the collection of genome sets associated with each unique and doubly-unique substring in the read  $r_j$ . If the intersection of these sets is empty, we discard the read since these substrings are “conflicting” (see below for a short discussion on these conflicts). If the intersection is non-empty, for each unique substring  $u_i$  and each doubly-unique substring  $d_i$  in  $r_j$ , we increase by 1 the associated counter  $c(u_i)$  and respectively  $c(d_i)$  we maintain. These counters indicate the number of “conflict-free” reads that include each such substring in  $\mathcal{Q}$ . As described below, these counters will be essential to actual query processing in Stage 2.

1. Identify the set  $U_j = \{u_{j,1}, \dots, u_{j,\ell}\} \subset \mathcal{U}$  of unique substrings in  $r_j$ ; similarly identify the set  $D_j = \{d_{j,1}, \dots, d_{j,\ell'}\} \subset \mathcal{D}$  of doubly-unique substrings in  $r_j$ .
2. (a) If  $U_j = D_j = \emptyset$  then discard  $r_j$ .

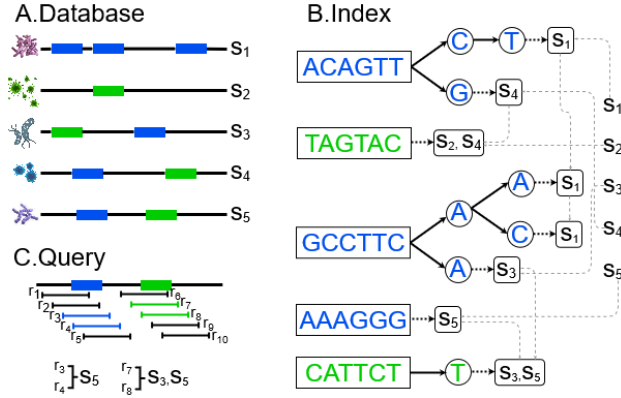


Figure 1: Overview of CAMMiQ's index structure. Strings in blue are unique substrings and those in green are doubly-unique substrings.

- (b) If  $U_j \neq \emptyset$  but it includes a pair of unique substrings  $u_{j,i}$  and  $u_{j,i'}$  which originate from different genomes, then discard  $r_j$ .
- (c) If  $U_j \neq \emptyset$  and all its unique substrings originate from the same genome  $s_k$ , however  $D_j$  includes a substring  $d_{j,i}$  which can not originate from  $s_k$ , then again discard  $r_j$ .
- (d) If  $U_j = \emptyset$  and the intersection of the set of genomes from where the substrings in  $D_j$  can originate is empty, then also discard  $r_j$ .
- (e) If on the other hand,
  - i.  $U_j \neq \emptyset$ , all its unique substrings originate from the same genome  $s_k$ , and each doubly-unique substring  $d_{j,i'} \in D_j$  can originate from  $s_k$ , or
  - ii.  $U_j = \emptyset$ , however the intersection between the genomes where the doubly-unique substrings of  $r_j$  can originate from is comprised of only one genome,  $s_k$ , or
  - iii.  $U_j = \emptyset$  and  $d_j$  is comprised of doubly-unique substrings that can only originate from the same pair of genomes  $s_k$  and  $s_{k'}$ , then
 then increase  $c(u_{j,i})$  by 1 for each  $u_{j,i} \in U_j$  and  $c(d_{j,i})$  by 1 for each  $d_{j,i} \in D_j$ .

The above counters are sufficient to compute the set  $\mathcal{A}_1$  as well as  $\mathcal{A}_3$ , the answer to our most general query type. For computing  $\mathcal{A}_2$ , CAMMiQ additionally maintains a counter  $d(s_k, s_{k'})$  for each pair of genomes  $s_k, s_{k'}$ , indicating the number of reads in  $\mathcal{Q}$  that can originate both from  $s_k$  and  $s_{k'}$ ; the value of this counter needs to be increased for each pair of involved genomes by 1 in case (iii) above.

The preprocessing stage described above eliminates those reads that include conflicting unique or doubly-unique substrings - conflicting in the sense that they are associated with different genomes. There are two main reasons for observing such conflicts: (i) sequencing errors, (ii) the presence of reads in the query from genomes that are not in  $\mathcal{S}$ . By eliminating these conflicting reads, we reduce the chances of mis-identifying the genomes from which they may originate.

In short, the read preprocessing stage produces two vectors  $\mathbf{c}_i^u = (c(u_{i,1}), \dots, c(u_{i,nu_i}))$  and  $\mathbf{c}_i^d = (c(d_{i,1}), \dots, c(d_{i,nd_i}))$  that indicate the number of (conflict-free) reads that include each unique and doubly-unique substring on each genome  $s_i$ . One can use these vectors to compute  $\mathcal{A}_1 = \{s_i : \sum_{l=1}^{nu_i} c(u_{i,l}) > 0\}$ . Additionally, through the use of the counter  $d(s_k, s_{k'})$  maintained for each pair of genomes  $s_k, s_{k'}$  one can compute  $\mathcal{A}_2 = \arg \min |\mathcal{A}' \subset \mathcal{S}|$  such that (i)  $s_i \in \mathcal{A}'$  if  $\sum_{l=1}^{nu_i} c(u_{i,l}) > 0$  and (ii)  $\exists s_i \in \mathcal{A}'$ , if  $d(s_k, s_{k'}) > 0$  then either  $i = k$  or  $i = k'$ . This is basically the solution to the hitting set problem we mentioned earlier, whose formulation as an integer linear program (ILP) is well known [56]. From this point on, our main focus will be computing  $\mathcal{A}_3$  (as well as the relative abundances of the genomes in  $\mathcal{A}_3$ ) for which we introduce an ILP formulation described below.



### 2.3 Query Processing Stage 2: ILP Formulation

CAMMiQ computes the list of genomes in the query as well as their abundances through an integer linear program (ILP) described below. Let  $\delta_i = 0/1$  be the indicator for the absence or presence of the genome  $s_i$  in  $\mathcal{Q}$ . The ILP formulation assigns a value to each  $\delta_i$  and also computes for each  $s_i$  its abundance  $p_i$  in the range  $[p_{\min}, p_{\max}]$ .

$$\begin{aligned} \text{Minimize} \quad & \sum_i \left( \frac{1}{nu_i} \sum_{l=1}^{nu_i} |c(u_{i,l}) - e(u_{i,l})| + \frac{1}{nd_i} \sum_{l=1}^{nd_i} |c(d_{i,l}) - e(d_{i,l})| \right) \\ \text{Subject to} \quad & e(u_{i,l}) = (L - |u_{i,l}|) \cdot p_i \cdot \frac{1}{L} \cdot (1 - \text{e}\hat{\text{r}})^{|u_{i,l}|} \quad \forall i, l, \text{ s.t. } 1 \leq l \leq nu_i \end{aligned} \quad (7)$$

$$e(d_{i,l}) = (L - |d_{i,l}|) \cdot (p_i + p_j) \cdot \frac{1}{L} \cdot (1 - \text{e}\hat{\text{r}})^{|d_{i,l}|} \quad \forall i, l \text{ s.t. } 1 \leq l \leq nd_i \quad (8)$$

$$p_i \leq \delta_i \cdot p_{\max} \quad \forall i \quad (9)$$

$$\delta_i = 0 \quad \forall i, \text{ s.t. } s_i \notin E(\mathcal{Q}) \quad (10)$$

$$p_i \geq \delta_i \cdot \min \left\{ L \cdot \sum_{l=1}^{nu_i} c(u_{i,l}) \cdot \frac{1}{nu_i^L}, L \cdot \sum_{l=1}^{nd_i} c(d_{i,l}) \cdot \frac{1}{nd_i^L} \right\} \cdot (1 - \epsilon) \quad \forall i, \text{ s.t. } s_i \in E(\mathcal{Q}) \quad (11)$$

$$\sum_i |s_i| \cdot p_i \leq n \cdot L \quad (12)$$

The objective of the ILP is to minimize the sum of absolute difference between the expected and the actual number of reads to cover a unique or doubly-unique substring. Since each genome may have different number of unique and doubly-unique substrings, this difference is normalized w.r.t.  $nu_i$  or  $nd_i$ . Constraints (7) and (8) define the expected number of reads to cover a particular unique substring  $u_{i,l}$  or doubly-unique substring  $d_{i,l}$  - here  $p_j$  is the abundance of genome  $s_j$  which also includes  $d_{i,l}$ . Here we denote by  $\text{e}\hat{\text{r}}$  the estimated read error (specifically substitution) rate per nucleotide, and denote by  $|w|$  the length of a substring  $w$ . Constraint (9) ensures that the abundance  $p_i$  of a genome is 0 if  $\delta_i = 0$ . Constraint (10) ensures that the solution to the above ILP excludes those genomes whose counters for unique and doubly-unique substrings add up to a value below a threshold - so as to reduce the size of the solution space. More specifically, given a threshold value  $\alpha$ , the constraint excludes those genomes  $s_i$  that are not in the set of genomes  $E(\mathcal{Q})$  whose counters for its unique substrings add up to a value above  $\alpha \cdot nu_i^L$ , and doubly unique substrings add up to a value above  $\alpha \cdot nd_i^L$ . More formally,  $E(\mathcal{Q}) = \{s_i \in \mathcal{S} \mid \sum_{l=1}^{nu_i} c(u_{i,l}) \geq \alpha \cdot nu_i^L\} \cap \{s_i \in \mathcal{S} \mid \sum_{l=1}^{nd_i} c(d_{i,l}) \geq \alpha \cdot nd_i^L\}$ . Constraint (11) enforces a lower bound on the coverage (and thus the abundance) of each genome  $s_i$  in the solution to the above ILP (namely, with  $\delta_i = 1$ ), which must match the coverage  $(L \cdot \sum_{l=1}^{nu_i} c(u_{i,l}) \cdot \frac{1}{nu_i^L}$  and  $\sum_{l=1}^{nd_i} c(d_{i,l}) \cdot \frac{1}{nd_i^L}$ ) resulting from the number of reads in  $\mathcal{Q}$  that include a unique and doubly-unique substring respectively, i.e., it must be at least  $(1 - \epsilon)$  times the smaller one above for a user defined  $\epsilon$ . Constraint (12) enforces an upper bound on the coverage (and thus the abundance) of each genome  $s_i$  in the solution to the above ILP, through making the sum (over each  $s_i$ ) of the number of reads produced on  $s_i$  based on  $p_i$  not exceed the total number of reads  $n$ . Collectively, the last two constraints ensure that the abundance  $p_i$  computed from the ILP matches what is (i.e., the coverage based on read counts) given by  $\mathcal{Q}$ . Note that the absolute values in the objective can be removed by introducing a new variable  $\gamma(u_{i,l}) \geq \max\{c(u_{i,l}) - e(u_{i,l}), e(u_{i,l}) - c(u_{i,l})\}$ .

### 2.4 When to Use Unique Substrings - the Error Free Case

We now provide a set of sufficient conditions to guarantee the (approximate) performance that can be obtained (with high probability) in metagenomic identification and quantification by the use of unique

substrings only. These conditions apply to CAMMiQ when  $c = 1$ , as well as CLARK, KrakenUniq and other similar approaches. In case these conditions are not met, it is advisable to use CAMMiQ with  $c \geq 2$ .

Suppose that we are given a query  $\mathcal{Q}$  composed of  $n$  error-free reads of length  $L$ , sampled independently and uniformly at random from a collection of genomes  $\mathcal{A} = \{s_1, \dots, s_a\}$  according to their abundances  $p_1, \dots, p_a$ . More specifically, suppose that our goal is to answer query  $\mathcal{Q}$  by computing  $\mathcal{A}_1$ , along with an estimate for the abundance value  $p_i$  for each  $s_i \in \mathcal{A}_1$ , calculated as the weighted number of reads assigned to  $s_i$  according to the procedure described in Section 2.2. Then, the L1 distance between the true abundance values and this estimate will not exceed a value determined by  $n$  (number of reads),  $a$ , and  $q_{\min}$ , the minimum (normalized) proportion of unique  $L$ -mers among these genomes. For a given failure probability  $\delta$  and an upper bound on L1 distance  $\epsilon$ , this translates into sufficient conditions on the values of  $n$ ,  $a$  and  $q_{\min}$  to ensure acceptable performance by the computational method in use.

**Theorem 1.** *Let  $\mathcal{Q} = \{r_1, \dots, r_n\}$  be a set of  $n$  error-free reads of length  $L$ , each sampled independently and uniformly at random from all positions on a genome  $s_i \in \mathcal{A} = \{s_1, \dots, s_a\}$ , where  $s_1, \dots, s_a$  is distributed according to their abundances  $p_1, \dots, p_a > 0$ . Let  $p'_i = \frac{p_i \cdot n_i^L}{\sum_{i'=1}^a p_{i'} \cdot n_{i'}^L}$  be the corresponding “unnormalized” abundance of  $p_i$  for  $i = 1, \dots, a$ , where  $n_i^L$  denotes the total number of  $L$ -mers on  $s_i$ . Let  $q_1, \dots, q_a > 0$  be the proportion of unique  $L$ -mers on  $s_1, \dots, s_a$  respectively;  $p_{\min} = \min\{p_1, \dots, p_a\}$ ;  $q_{\min} = \min\{q_1, \dots, q_a\}$ . Then,*

- (i) *With probability at least  $1 - \delta$ , each  $s_i$  can be identified through querying  $\mathcal{Q}$  if  $n \geq \frac{2(a+1)+\ln(1/\delta)}{(p_{\min} q_{\min})^2}$*
- (ii) *With probability at least  $1 - \delta$ , the L1 distance between the predicted abundances  $\hat{p}_1, \dots, \hat{p}_a$  by setting  $\hat{p}_i = \frac{c_i/q_i}{n}$  and the true (unnormalized) abundances  $p'_1, \dots, p'_a$  is at most  $\epsilon$  if  $n \geq \frac{2(a+1)+\ln(1/\delta)}{(\epsilon q_{\min})^2}$ .*
- (iii) *Given  $n$  such reads in a query, with probability at least  $1 - \delta$ , the L1 distance between the predicted abundances  $\hat{p}_1, \dots, \hat{p}_a$  by setting  $\hat{p}_i = \frac{c_i/q_i}{n}$  and the true (unnormalized) abundances  $p'_1, \dots, p'_a$  is bounded by  $\sqrt{\frac{2[\ln(1/\delta)+(a+1)]}{n q_{\min}^2}}$ .*

where  $c_i$  denotes the number of reads assigned to  $s_i$ .

See Supplementary Methods for a proof.

### 3 Results

In order to evaluate the overall running time, memory utilization and accuracy of CAMMiQ, we have established three datasets to be indexed, all based on NCBI’s RefSeq database [57]. The first dataset aims to evaluate the species level performance of CAMMiQ and thus includes one representative genome from each bacterial species. The next dataset aims to evaluate CAMMiQ’s strain level performance by including all strain level genome data from RefSeq. Since RefSeq database is biased towards a handful of genera (e.g. *Salmonella* genus is represented by more than 300 strains) we also generated a subspecies level dataset where each species is represented by only a handful of strains as described in detail below.

On the three index datasets, we evaluated CAMMiQ on both simulated and real bacterial read sets. The simulated query sets aimed to measure the comparative performance of CAMMiQ against some of the best used methods available. The real query sets aimed to evaluate the accuracy of CAMMiQ on metatranscriptomic reads from single human immune cells deliberately infected with two distinct strains of the intracellular bacterium *Salmonella enterica* [58].

For each dataset to be indexed and each corresponding query set we compared CAMMiQ’s performance with some of the best computational methods available. As will be demonstrated, CAMMiQ’s performance is superior to all alternatives in almost all scenarios we tested.

Below, in Section 3.1 we give a detailed description of the index data sets, simulated and real query sets we used, as well as the alternative computational methods we tested to benchmark CAMMiQ’s performance. Then in Section 3.2, we demonstrate the advantage offered by CAMMiQ’s (i) utilization of doubly-unique substrings - in addition to unique substrings, and (ii) consideration of shortest such substrings instead of fixed length  $k$ -mers as per CLARK. In Sections 3.3 and 3.4 we demonstrate CAMMiQ’s comparative accuracy against

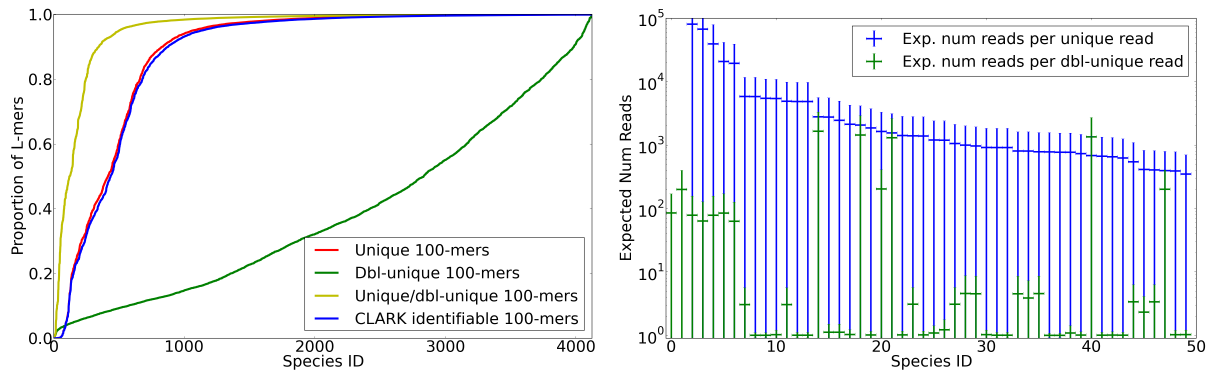


Figure 2: (A) Proportion of  $L$ -mers (for  $L = 100$ ) that include a unique substring (plotted in red), a doubly-unique substring (plotted in green), or either a unique or a doubly-unique substring (plotted in yellow) in the large dataset of 4122 bacterial genomes from NCBI RefSeq database; these  $L$ -mers, when presented as reads are utilizable by CAMMIQ (when  $L_{\max} = 100$ ). Also included are the proportion of  $L$ -mers that include a unique  $k$ -mer ( $k = 30$ ) and thus are utilized by CLARK. For each plot, the genomes are independently sorted with respect to the corresponding proportions in ascending order. (B) The expected number of reads (of length  $L = 100$ ) needed to capture one read containing a unique substring (in blue) as well as one read containing a doubly-unique substring (in green) in the 50 genomes with the lowest proportion of unique  $L$ -mers. Also included are the range for each value within the corresponding standard deviation. Note that for this figure  $L_{\max} = 50$ .

alternative metagenomic analysis methods on simulated query sets through the use of our species-level and strain-level datasets. CAMMIQ’s comparative performance on computational resources is then demonstrated in Section 3.5. Finally in Section 3.6 we demonstrate CAMMIQ’s performance on real metatranscriptomic query sets through its use of our the subspecies-level dataset. The results of CAMMIQ in this setup was compared against that of GATK PathSeq which was used for the same purpose earlier [14], as well as BLASTN [9].

### 3.1 The Experimental Setup: Index Datasets, Queries, Benchmarked Methods, Hardware

Our “species-level” index dataset consists of all complete bacterial genomes in NCBI’s RefSeq [57] Database (downloaded on 06/16/2019). We (randomly) selected one representative genome per species out of 13737 reference genomes representing 4122 distinct species. This resulted in a total of  $m = 4122$  genomes with a total length of  $M = 3.4 * 10^{10}$  (including the reverse complement contigs). On this index dataset we used simulated queries only - as will be described below. Since our goal here is to measure CAMMIQ’s relative performance against available methods, we also benchmarked state of the art  $k$ -mer based and marker-gene based metagenomics profiling tools on this index dataset. Specifically we used Kraken2 [59] (the latest version of Kraken [12]), KrakenUniq [44], CLARK [26] and MetaPhlan2 [60]; these four tools provide very similar functionality to CAMMIQ such as read level classification and abundance estimation.

Our “strain-level” index dataset is smaller: it is restricted to 614 Human Gut related genomes (according to [61]) and is designed to evaluate CAMMIQ’s strain level identification and quantification performance against the available tools mentioned above. Again, we used simulated query sets for this index dataset. Note that on both the strain-level and species-level index datasets, we evaluated CAMMIQ on the most general  $\mathcal{A}_3$  type queries.

Perhaps our most interesting results are on the “subspecies-level” index dataset, which consists of 3395 selected genomes from the 13737 complete bacterial genomes; this dataset was primarily designed to evaluate CAMMIQ’s accuracy on query sets involving real metatranscriptomic samples obtained from 262 single human immune cells [58], each exposed to or infected with one of the two distinct *Salmonella* strains. In addition there are 80 uninfected cells used as negative controls. The reads from each cell forms a natural query set for this index dataset. A recent study [14] used the GATK PathSeq tool [10] in order to validate the *Salmonella*

Query Set	Index Dataset	Read Length ( $L$ )	Num. Reads ( $n$ )	Err Rate
Least-20-uniform-e.f.		100	4.8M	0
Least-20-uniform		100	4.8M	0.01
Least-quantifiable-20-uniform-e.f.		100	5.0M	0
Least-quantifiable-20-uniform		100	5.0M	0.01
Least-20-genera-uniform-e.f.	Species	100	4.0M	0
Least-20-genera-uniform	level	100	4.0M	0.01
Least-20-genera-lognormal	dataset	100	4.0M	0.01
Random-20-uniform		100	4.4M	0.01
Random-20-lognormal		100	5.0M	0.01
Random-20-lognormal-a.g.		100	1.1M	0.01
Random-100-uniform		100	21.5M	0.006
HumanGut-least-25	Strain	100	2.0M	0.01
HumanGut-random-100-1	level	100	8.0M	0.01
HumanGut-random-100-2	dataset	125	8.0M	0.01
HumanGut-all		100	20.0M	0.01
Filtered-scRNA-seq	Subspecies level dataset	66.4(13.6)	8.5M	Unknown

Table 1: Synthetic and real bacterial read sets used to benchmark CAMMiQ’s performance against the best available metagenomic classification and abundance estimation tools. The top collection of queries involve genomes from our species-level index dataset consisting of 4122 distinct bacterial species. The middle collection were sampled from our strain-level index dataset of 617 gastrointestinal associated bacteria (possibly incompletely assembled) provided in [61]. The bottom collection were real single cell RNA-seq reads sequenced from 342 immune cells infected with *Salmonella enterica* [58]; the corresponding subspecies-level index dataset consists a selection of 3395 bacterial genomes from 2753 distinct species.

strains associated with each cell with limited success. Unlike the tools benchmarked above, GATK PathSeq is alignment based; as a consequence it is slower than the above alternatives but is expected to be more accurate. Interestingly CAMMiQ’s ability to distinguish cells exposed to or infected with specific strains of *Salmonella* was better than PathSeq’s ability to do the same - with the added bonus that it is much faster. Note that on metatranscriptomic query sets we used query types  $\mathcal{A}_1$  and  $\mathcal{A}_2$  rather than the most general query type  $\mathcal{A}_3$ .

In order to systematically compare the performance of CAMMiQ and other metagenomic profiling tools, we generated several simulated metagenomes as query sets, summarized in Table 1. The upper part of the table corresponds to simulated data from our species-level dataset and the middle part corresponds to our strain-level dataset.

The first set of simulated metagenomes aim to assess how well CAMMiQ identifies species in a query. For that we simulated a metagenome consisting of the 20 genomes that have the lowest number of unique  $L$ -mers in our species-level dataset. Each genome in the mixture was simulated to have similar read coverage. The very first query we generated from this mixture (denoted Least-20-uniform-e.f.) had no read errors. The second query (denoted Least-20-uniform) had i.i.d. substitution errors occurring at a rate of 1%. Note that the 20 genomes we used in this mixture are intrinsically difficult to be identified by CLARK and other tools we compared. However since these genomes have many doubly-unique  $L$ -mers - which are sometimes shared with multiple other genomes, they could be identified by CAMMiQ (see Supplementary Methods and Figure 5 for a more detailed explanation).

The second set of metagenomes we simulated aim to assess the species-level quantification performance of CAMMiQ. This simulated metagenome consisted of the 20 genomes which are among the 50 genomes in our species-level dataset with the lowest proportion of unique  $L$ -mers, but had the highest proportion of doubly-unique  $L$ -mers, making them somewhat easier to identify in comparison to the simulated metagenome above, but difficult to quantify by tools other than CAMMiQ. We again generated two simulated read collections

from this mixture, one with no sequencing errors (denoted Least-quantifiable-20-uniform-e.f.) and another with i.i.d. substitution rate of 1% (denoted Least-quantifiable-20-uniform). The ability of CAMMiQ’s ILP formulation to simultaneously determine the presence and abundance of genomes in these queries help it outperform the alternatives. (See Supplementary Methods and Figure 5 for a more detailed explanation.)

Even though RefSeq identifies each genome in our species-level dataset to represent a distinct species, a few of them have “unclassified” lineages at the species level. (Some of the genomes in the above queries are among them; see Figure 5 in the Supplementary Methods.) For example, *Rhizobium sp. N1314* with Taxonomy ID: 1703961 has Rank: species; however its Lineage is noted as *unclassified Rhizobium*. Because of this ambiguity, we generated a third set of metagenomes, again consisting 20 of those 50 genomes with the lowest proportion of unique  $L$ -mers, this time making sure that each of these 20 genomes represent a distinct genus. We generated 3 queries from this set of genomes. The first one had uniform coverage and had no sequencing errors (denoted Least-20-genera-uniform-e.f.). The second and third both had i.i.d. substitutions at a rate of 1%; the second had a uniform read coverage (denoted Least-20-genera-uniform), while the third had lognormal distribution (denoted Least-20-genera-lognormal).

In addition to the above particularly challenging queries, we simulated a number of additional read collections from 20 to 100 randomly chosen genomes from our species-level dataset. Unlike the above queries, all these read collections had i.i.d. substitution errors; the first three queries at a rate of 1% and the last query at a rate of 0.6%. Among them, the first simulated query (denoted Random-20-uniform) included reads from 20 genomes, each with similar read coverage. The second (denoted Random-20-lognormal) again included reads from 20 genomes, this time with coverages obeying a log-normal distribution. The third (denoted Random-20-lognormal-a.g.) included reads from 20 genomes, again with log-normal coverage distribution; what makes this query unique is that 10% of the reads were from an additional genome (denoted in the dataset name as a.g.) not included in our species-level dataset and thus is not part of CAMMiQ’s index. The fourth (denoted Random-100-uniform) included reads from 100 randomly chosen genomes from our species-level dataset, all with similar coverage.

In order to assess CAMMiQ’s performance in strain level identification and quantification, we simulated queries involving genomes from a database of 614 strains of human gastrointestinal bacteria [61]<sup>3</sup> from 409 species. We again simulated multiple queries, the first one involving reads from 25 strains with the smallest number of unique  $L$ -mers (denoted HumanGut-least-25), next two involving reads from randomly selected 100 strains, the first with  $L = 100$  as per the remainder of the queries (denoted HumanGut-random-100-1), and the second with  $L = 125$  (denoted HumanGut-random-100-2), and the final involving reads sampled from 409 randomly picked strain level genomes (denoted HumanGut-all), each representing a distinct species in the dataset. Note that none of these queries included more than one strain per species since two distinct strains from a species are not likely to be simultaneously present in a metagenomic sample.

As mentioned earlier, the index we built to respond to these queries consisted of all the 614 strain level genomes described above. The majority of these genomes are not complete and is comprised of multiple contigs; we filtered out any contig with length  $< 10$ KB and built the index on the remaining contigs. This resulted in seven strains without a single unique 100-mer and one strain without a single unique or doubly unique 100-mer. This last genome of *Bacillus andreae* was excluded from our queries since it contains no indexable substring.

In our final experiment, we applied CAMMiQ to a “gold-standard” query set consisting of immune cells infected *ex vivo* with an *intracellular* bacteria *Salmonella enterica* and subsequently sequenced using single-cell RNAseq (scRNA-seq) [58] to validate its feasibility of identifying microbial reads from real datasets. Specifically, this query, denoted as Filtered-scRNA-seq and also summarized in Table 1, consists of 342 monocyte-derived dendritic cells (moDCs) infected with either the *D23580* (STM-D23580) or the *LT2* (STM-LT2) strain of *Salmonella enterica* and sequenced using Smart-seq2 platform. The corresponding index we built to respond to this query consisted of the sparsified set of unique and doubly unique substring (with  $L = 75$ ,  $L_{\min} = 26$  and  $L_{\max} = 50$ ) from the subspecies-level dataset of 3395 selected complete bacterial genomes in NCBI’s RefSeq Database. To create this subspecies-level dataset, we first identified 2753 out of the total 4122 species to which at least 10 reads were mapped using GATK PathSeq [10] - and then we subsampled the 4325 (out of the total 13737) genomes from these 2753 species by only keeping one representative for each child of a species level taxonomic ID in NCBI’s Taxonomy Database; The only

<sup>3</sup>The complete set of genomes in this database is 617 but only 614 can be downloaded from RefSeq.



exception is we replaced the sampled genome of *Salmonella enterica subsp. enterica* (Taxonomy ID: 59201) with the genome of the two related strains *D23580* (Taxonomy ID: 568708) and *LT2* (Taxonomy ID: 99287), which results in the final 3395 subspecies level representatives. Note that as a final step (which is in contrast to the index used in the above simulated queries), we removed all plasmids contained in these genomes, since the plasmids in many genomes downloaded from RefSeq were missing, which may cause potential false positive “unique” or “doubly unique” substrings.

All of our experiments were run on a Linux server equipped with 40 Intel Xeon E7-8891 2.80 GHz processors, with 2.5 TB of physical memory and 30 TB of disk space. The ILP solver used by CAMMIQ is IBM ILOG CPLEX 12.9.0.

### 3.2 Comparative Utility of Variable-Length and Doubly-Unique Substrings

In this section we demonstrate the theoretical advantage offered by CAMMIQ in comparison to other tools, due to its utilization of not only unique  $k$ -mers (as per CLARK), but substrings of any length, which could be either unique or doubly-unique - within the species-level dataset we constructed. For that we compared the proportion of  $L$ -mers from each genome  $s_i$  in this dataset (for read length  $L = 100$ ) that are unique or doubly-unique (and thus can be utilized by CAMMIQ) with the proportion of  $L$ -mers that include a unique  $k$ -mer (that can be utilized by CLARK) for  $k = 30$ .

Figure 2A summarizes our findings: on the horizontal axis, the genomes are sorted with respect to the proportion of unique and doubly-unique  $L$ -mers they have; the vertical axis depicts this proportionality (from 0.0 to 1.0). The four plots we have are for the proportion of unique  $L$ -mers, doubly-unique  $L$ -mers, the combination of unique and doubly-unique  $L$ -mers (all utilized by CAMMIQ), and the  $L$ -mers that include a unique  $k$ -mer (utilized by CLARK). As can be seen, roughly three quarters of all genomes in this dataset (and thus many of the bacterial species) are easily distinguishable since a large fraction of their  $L$ -mers include a unique  $k$ -mer. However, roughly a quarter of the genomes in this dataset can benefit from the consideration of doubly-unique substrings, especially when their abundances are low. In particular, 66 of these 4122 genomes/species have extremely low proportion of ( $\leq 1\%$ ) unique 100-mers. In fact, the species *Francisella sp. MA06-7296* does not have a single unique 100-mer and the species *Rhizobium sp. N6212* does not have any 100-mer that include a unique 30-mer (in fact any substring of length  $\leq L_{\max} = 50$ ). These two species cannot be identified by CLARK in any microbial mixture, independent of the abundance values.

Figure 2B depicts the inverse proportionality of doubly-unique  $L$ -mers in comparison to unique  $L$ -mers among 50 genomes that have the lowest proportion of unique  $L$ -mers - for  $L = 100$ . The inverse-proportionality of unique or doubly-unique  $L$ -mers for a genome corresponds to the number of reads to be sampled (on average) from that genome to guarantee that the sample includes one read that is guaranteed to be assigned to the correct genome. In the absence of read errors, this guarantees the correct identification of the corresponding genome in the query. Note that, in half of these 50 genomes, almost all  $L$ -mers are doubly-unique. This implies that any query involving one or more of these genomes would unlikely to be resolved with CLARK. Yet this query could be handled by CAMMIQ since if a read contains a unique or doubly-unique substring, then it can be correctly assigned to the corresponding species as described in Section 2.2.

Note that 3296 of the 4122 species in our species-level dataset have  $\geq 90\%$  of their 100-mers as unique. Only a few of these unique substrings do not include a unique 30-mer and thus will be missed by CLARK. This implies that from the accuracy point of view on these genomes, CLARK’s use of  $k$ -mers instead of the shortest substrings does not put it at a disadvantage on these genomes - when  $k = 30$  and  $L = 100$ . However, as it will become clear later (see Table 5) CAMMIQ’s use of shortest substrings, combined with its subsampling strategy gives it an advantage over CLARK (as well as KrakenUniq) with respect to the size of the index structure - despite the fact that CAMMIQ needs to index not only unique but also doubly-unique substrings from each genome.

### 3.3 Classification and Quantification Performance at the Species Level

We tested CAMMIQ on read collections sampled from the species-level dataset and compared its practical performance with the select metagenomic profiling tools. Perhaps the most widely-used performance measures to benchmark metagenomic classifiers are the proportion of reads correctly assigned to a genome among

Performance Measure	Query Set	CAMMiQ	Kraken2	KrakenUniq	CLARK	MetaPhlAn2
A. Precision	Least-20-uniform-e.f.	<b>1.0</b>	0.014	0.952	<b>1.0</b>	-
	Least-20-uniform	<b>0.974</b>	0.001	0.009	0.006	-
	Least-quantifiable-20-uniform-e.f.	<b>1.0</b>	0.043	0.936	<b>1.0</b>	-
	Least-quantifiable-20-uniform	<b>0.989</b>	0.012	0.086	0.064	-
	Least-20-genera-uniform-e.f.	<b>1.0</b>	0.027	0.978	<b>1.0</b>	-
	Least-20-genera-uniform	<b>0.993</b>	0.005	0.047	0.028	-
	Least-20-genera-lognormal	<b>0.993</b>	0.005	0.049	0.033	-
	Random-20-uniform	<b>0.997</b>	0.783	0.981	0.994	-
	Random-20-lognormal	<b>0.998</b>	0.933	0.991	0.995	-
	Random-20-lognormal-a.g.*	<b>0.998</b>	0.900	<b>0.998</b>	0.997	-
Random-100-uniform	<b>0.998</b>	0.968	0.965	0.997	-	
B. Num. Assigned Reads	Least-20-uniform-e.f.	<b>1.72M</b>	25226	664	675	76483
	Least-20-uniform	<b>1.57M</b>	220721	62189	94632	70605
	Least-quantifiable-20-uniform-e.f.	<b>3.18M</b>	79408	6307	5844	120199
	Least-quantifiable-20-uniform	<b>2.81M</b>	246877	62617	81575	115959
	Least-20-genera-uniform-e.f.	<b>2.89M</b>	33523	1844	1763	67625
	Least-20-genera-uniform	<b>2.58M</b>	170506	36286	57764	64204
	Least-20-genera-lognormal	<b>2.53M</b>	183709	41681	59679	73433
	Random-20-uniform	3.84M	3.33M	<b>3.88M</b>	3.81M	80633
	Random-20-lognormal	<b>4.48M</b>	3.84M	4.43M	4.39M	19484
	Random-20-lognormal-a.g.*	0.91M	0.93M	<b>0.97M</b>	<b>0.97M</b>	28328
Random-100-uniform	<b>19.5M</b>	17.0M	19.4M	18.7M	0.4M	
C. Num. Correctly Identified Genomes	Least-20-uniform-e.f.	20/20	16	17	18	13
	Least-20-uniform	20/28	16	17	18	13
	Least-quantifiable-20-uniform-e.f.	20/20	19	19	20	18
	Least-quantifiable-20-uniform	20/27	19	19	20	18
	Least-20-genera-uniform-e.f.	20/20	18	18	18	12
	Least-20-genera-uniform	20/24	18	18	18	12
	Least-20-genera-lognormal	20/33	17	17	18	12
	Random-20-uniform	20/20	18	20	20	11
	Random-20-lognormal	20/21	20	20	20	9
	Random-20-lognormal-a.g.*	20/20	19	20	20	13
Random-100-uniform	100/100	100	100	100	76	
D. L1 Err.	Least-20-uniform-e.f.	<b>0.0790</b>	-	-	0.8846	0.8171
	Least-20-uniform	<b>0.0929</b>	-	-	0.9124	0.8180
	Least-quantifiable-20-uniform-e.f.	<b>0.0375</b>	-	-	0.5774	0.5722
	Least-quantifiable-20-uniform	<b>0.0278</b>	-	-	0.5659	0.5956
	Least-20-genera-uniform-e.f.	<b>0.0626</b>	-	-	0.7156	0.9153
	Least-20-genera-uniform	<b>0.0591</b>	-	-	0.7067	0.9533
	Least-20-genera-lognormal	<b>0.0439</b>	-	-	0.5701	1.0083
	Random-20-uniform	<b>0.0113</b>	0.4139	0.1446	0.2042	0.8257
	Random-20-lognormal	<b>0.0038</b>	0.2843	0.1217	0.1496	1.7367
	Random-20-lognormal-a.g.*	0.1262	0.2412	<b>0.1173</b>	0.1861	0.5578
Random-100-uniform	<b>0.0096</b>	0.2364	0.1327	0.2063	0.7801	

Table 2: Performance evaluation of CAMMiQ, Kraken2, KrakenUniq, CLARK and MetaPhlAn2 on queries from the species level dataset. Precision: the proportion of reads correctly assigned to a genome among the set of reads assigned to some genome (correctly or incorrectly). Number of assigned reads: the total number of reads assigned to some genome. Number of correctly identified genomes: for CAMMiQ we report both the number of correctly identified genomes (true positives) and the total number of genomes returned by its ILP formulation (false positive); for Kraken2, KrakenUniq, CLARK and MetaPhlAn2 we only report the number of correctly identified genomes (true positives). Note that we consider MetaPhlAn2 to have a correct identification even if it reports the genus that this genome belongs to. L1 error: the L1 distance between the true relative abundance values (between 0 and 1) and the predicted abundance values for each genome in the query (i.e. positives). We made an exception for MetaPhlAn2, where we measured the genus level L1 distance. Note that we converted the true abundance values reported by Kraken2, KrakenUniq and CLARK by dividing the predicted abundance value for each genome by its length and then normalizing these values by the total abundance value of all genomes. In each of the four measures, the best performing tool’s results are highlighted.

\*: 10% reads in the query Random-20-lognormal-a.g. are from a genome not in the index; any assignment of these reads are necessarily incorrect by all tools except MetaPhlAn2 - which uses its own index, that happens to include this genome.

(i) the set of reads assigned to some genome, i.e. *precision*, and (ii) the full set of reads in the query, i.e. *sensitivity* [26]. In Table 2A we report the precision values for all tools we benchmarked with the exception of MetaPhlAn2, since it does not report individual read assignments. In Table 2B, instead of reporting sensitivity, we report the total number of reads assigned to a genome, since the main goal of Kraken2, KrakenUniq and CLARK is to classify as many reads as possible correctly. The number of reads assigned to a genome (part B of the Table) multiplied by precision (part A of the table) indeed gives this number, i.e. that of correctly classified reads. Note that any read assigned to a taxonomic level higher than (and not including) species level by Kraken2 or KrakenUniq are considered to be not assigned. (In a way, this increases their reported precision but decreases their number of assigned reads.) Also note that CAMMIQ’s primary goal is not to classify each read but rather identify and quantify genomes in a query. Nonetheless, we still report its intermediate output as follows. CAMMIQ considers certain reads as conflicting; here we consider them as not assigned. It assigns certain reads to a single genome; we consider each such read assigned, and if the assignment is correct, also correctly assigned. CAMMIQ then assigns each remaining read ambiguously to two potential genomes.<sup>4</sup> We consider this read assigned, and in case one of these two genomes is correct, also correctly assigned. Finally since MetaPhlAn2 has an index based on a predetermined database, it is not easy to evaluate this tool’s precision, since the IDs of the genomes indexed by the other four tools do not always correspond to MetaPhlAn2’s genome IDs. As a result we only report its number of assigned reads but not its precision.

We benchmarked CAMMIQ using its default parameter settings of  $L_{\min} = 26$  and  $L_{\max} = 50$ , against Kraken2, KrakenUniq and CLARK, with all three using  $k$ -mer length of 26. All four of these tools used the same genomes for establishing their index. As can be seen in Table 2, CAMMIQ demonstrated the best precision for read classification in all of the 11 simulated query sets. With respect to the total number of assigned reads (correctly or incorrectly) on the first 7 queries, i.e. those involving the 20 genomes with the least number of unique  $L$ -mers (Least-20-uniform-e.f. and Least-20-uniform), those that are the least quantifiable (Least-quantifiable-20-uniform-e.f. and Least-quantifiable-20-uniform), and those are composed of genomes each from a distinct genus (Least-20-genera-uniform-e.f., Least-20-genera-uniform and Least-20-genera-lognormal), CAMMIQ improves over the alternatives not only in terms of precision but also the number of reads assigned, sometimes by an order of magnitude or more. The only exceptions are on those “hypothetical” error-free queries, on which not only CAMMIQ and but also CLARK achieves 100% precision. For the above 7 queries Kraken2 has the lowest precision; KrakenUniq has better precision but only to a degree. Furthermore, the number of assigned reads by KrakenUniq is typically the lowest.

On the next four queries, which are easier to identify and quantify, CAMMIQ’s performance is still the best overall. Its precision is the best for all four of these queries while the number of assigned reads are slightly worse than KrakenUniq in two queries. This is likely due to the fact that KrakenUniq produces some incorrect assignments, especially for the Random-20-uniform query on which KrakenUniq has a lower precision. The precision of CAMMIQ is identical to KrakenUniq on Random-20-lognormal-a.g. query, where 10% of the reads are sampled from an additional genome not indexed. This is despite that the number of assigned reads from this query is higher for KrakenUniq and CLARK, demonstrating that the introduction of un-indexed species impacts the performance of CAMMIQ similarly to the other tools.

We next evaluated the number of correctly identified genomes (for MetaPhlAn2, correctly identified species) in each query, as well as the L1 distance between the true abundance profile and the predicted abundance profile by all five tools on the 11 queries involving our species level dataset. The results can be found in Table 2, parts C and D and Figure 3. Note that CAMMIQ and MetaPhlAn2 automatically incorporate a normalization with respect to genome lengths, while Kraken2, KrakenUniq and CLARK simply report the number of reads assigned to each taxonomy rank (for our queries, species) as their abundance profile. In order to compute L1 distances correctly, we converted the true abundance profiles of Kraken2, KrakenUniq and CLARK by dividing the predicted abundance value of each genome by its length and then normalizing these values by the total abundance value of all genomes.

As can be seen in Table 2 panels C and D, as well as Figure 3, CAMMIQ clearly offers the best performance in both identification and quantification. It correctly identified all genomes present in each one of the 11 queries and was not impacted by the additional genome we introduced in the Random-20-lognormal-a.g. query. As importantly, it only returned very few false positive genomes for the most challenging Least-20-

<sup>4</sup>This happens if the read includes one or more doubly unique substrings from the same pair of genomes but no unique substring.

Query Set	$L_{\min}$	Precision	Num. Assigned Reads	Num. Identified Species		L1 Err.	
				$\alpha = 0.001$	$\alpha = 0.0001$	$\alpha = 0.001$	$\alpha = 0.0001$
Least-20-uniform	21	0.887	1.69M	20/21	31/32	0.0835	0.0836
	26	0.974	1.57M	20/21	28/29	0.0929	0.0929
	31	0.980	1.49M	20/20	22/23	0.1064	0.1060
Least-quantifiable-20-uniform	21	0.967	2.85M	20/20	29/31	0.0423	0.0400
	26	0.989	2.81M	20/20	27/29	0.0294	0.0278
	31	0.992	2.70M	20/20	24/26	0.0344	0.0326
Least-20-genera-uniform	21	0.974	2.61M	20/20	25/26	0.0706	0.0715
	26	0.993	2.58M	20/20	24/26	0.0585	0.0591
	31	0.995	2.48M	20/20	23/24	0.0688	0.0683
Least-20-genera-lognormal	21	0.974	2.56M	19/19	34/35	0.0394	0.0418
	26	0.993	2.53M	19/19	33/34	0.0416	0.0439
	31	0.995	2.43M	19/19	32/33	0.0355	0.0374
Random-20-uniform	21	0.993	3.83M	20/20	20/20	0.0220	0.0220
	26	0.997	3.84M	20/20	20/20	0.0113	0.0113
	31	0.998	3.69M	20/20	20/20	0.0294	0.0294
Random-20-lognormal	21	0.996	4.50M	20/20	21/21	0.0432	0.0431
	26	0.998	4.48M	20/20	21/21	0.0039	0.0038
	31	0.998	4.32M	20/20	21/21	0.0062	0.0058
Random-20-lognormal-a.g.*	21	0.989	0.92M	17/17	20/20	0.1492	0.1268
	26	0.998	0.91M	16/16	20/20	0.1631	0.1262
	31	0.999	0.87M	16/16	20/20	0.1658	0.1298
Random-100-uniform	21	0.996	19.5M	100/100	100/100	0.0176	0.0176
	26	0.998	19.5M	100/100	100/100	0.0096	0.0096
	31	0.999	19.0M	100/100	100/100	0.0104	0.0104

Table 3: Performance of CAMMiQ as a function of minimum unique/doubly-unique substring length  $L_{\min} = 21, 26, 31$ , and minimum relative read count threshold  $\alpha = 0.001, 0.0001$  to report a genome. Precision: the proportion of reads correctly assigned to a genome among the set of reads assigned to some genome, correctly or incorrectly. Number of assigned reads: the total number of reads assigned to some genome. Number of identified genomes: the number of genomes returned by CAMMiQ’s ILP formulation v.s. the number of genomes that have sufficient read assignments (determined by  $\alpha$ ). L1 error: the L1 distance between the true relative abundance values (between 0 and 1) and the predicted abundance values for each genome in the corresponding query.

\*: 10% reads in the query Random-20-lognormal-a.g. are from a genome not in the index; any assignment of such a read to a genome is necessarily incorrect.

uniform, Least-quantifiable-20-uniform and Least-20-genera-uniform as well as Least-20-genera-lognormal queries, and at most one false positive genome for the remaining 7 queries. Other tools had varying levels of false negatives in these 11 queries. Among them CLARK offered the best false negative performance: it only reported false negatives in the queries involving a genome without a single unique  $k$ -mer. As can be expected, MetaPhlAn2 had the worst performance with respect to false negatives, very likely due to the incompleteness of its marker gene list<sup>5</sup>. This also led to a larger L1 distance than the other tools, even for the relatively easy queries. Kraken2 and KrakenUniq were also prone to have a few false negatives, though less than MetaPhlAn2. Furthermore their predicted abundances are smaller than the true abundance values, because they may assign reads to higher taxonomic ranks than the species level; see Figure 3. This is especially evident in the in the first 7 (difficult) queries: even though Kraken2 and KrakenUniq identified the majority of the genomes correctly, the abundance values they reported on these genomes were all close to 0 and thus the L1 distances turned out to be very close to 1 (see Figure 3). As a result, these values are not reported in Table 2D. Note that we do not report the number of false positive genomes returned by each tool on these 11 queries. This is primarily because of the fact that CLARK, Kraken2 and KrakenUniq are not designed to identify genomes in a query mixture and thus do not “care” whether the false positive assignments are distributed across many genomes (which would result in a large number of false positives) or are concentrated in a few genomes (resulting in a few false positives). MetaPhlAn2 aims to identify genomes however it can do so in any taxonomic level. As such, all these four tools return a very large number of false positives, especially for the first 7 queries, but since we felt that this could be unfair due to the above reasons, we decided to report only the false positives reported by CAMMIQ, noting that its genome identification performance is significantly better than the other tools. We ignore the false positive calls also in the L1 measure since it is calculated only on the true positive genomes. This explains the single query and measure for which CAMMIQ seems to have performed worse than an alternative, namely KrakenUniq: on Random-20-lognormal-a.g. KrakenUniq’s L1 distance is slightly better than CAMMIQ. We remind that in this query, 10% of the reads are sampled from a genome not in the index. Since these reads can not be assigned through any means, the relative abundances of the other genomes will be overestimated by  $\sim 10\%$ , provided that there are no false positive genomes as per CAMMIQ. However, KrakenUniq’s false positive calls (there are several) reduces its relative abundance estimates for the true positive genomes, and gives a seemingly better L1 distance. This can be observed in Figure 3, which depicts for each query, relative abundance estimate for each genome by each one of the tools we benchmarked. As can be clearly seen, CAMMIQ’s individual abundance estimates are just on the mark for each one of the genomes even for the most difficult queries.

We finally evaluated the impact of two important parameters for CAMMIQ:  $\alpha$ , the minimum relative read count threshold for reporting a genome, and  $L_{\min}$ , the minimum unique or double-unique substring length (values larger than the default value of 50 for  $L_{\max}$  did not have a big impact and thus are not reported here). In Table 3 we report the results for each possible combination of  $L_{\min} = 21, 26, 31$  and  $\alpha = 0.001, 0.0001$ , on 8 of the 11 queries, omitting the 3 error free queries (on which the impact on precision is minimal). As we increase  $L_{\min}$ , CAMMIQ’s precision improves, however its read assignment performance deteriorates. Interestingly, its predicted abundance values did not change much with increasing  $L_{\min}$ . As a result we set the default  $L_{\min}$  to 26 in CAMMIQ. On the other hand, increasing the value of  $\alpha$ , decreased the number of false positives in CAMMIQ’s output, particularly in the most difficult queries. However, as a result of this, for the queries Least-20-genera-lognormal and Random-20-lognormal-a.g., those genomes with low abundance values were disregarded by CAMMIQ, leading to false negatives. CAMMIQ allows the user to set the parameter  $\alpha$  with prior knowledge on the reads to be queried (e.g., the expected read coverage or the number of genomes in the query); we set its default value to 0.0001.

### 3.4 Quantification Performance at the Strain Level

We finally evaluated CAMMIQ’s performance on queries composed from our strain-level dataset that consists of 614 Human Gut related genomes (strains) from 409 species [61] as described in Section 3.1. Three of the tools we evaluated, namely CLARK, Kraken2 and KrakenUniq, simply aim to perform read classification, which we evaluated in Section 3.3. They are not designed to identify or quantify genomes, especially at the strain level. On the other hand MetaPhlAn2’s index have strain level information and through that it attempts to identify and quantify strains. As a result we report on the performances of both CAMMIQ and MetaPhlAn2

<sup>5</sup>Here we used the latest set of marker genes `mpa_v20_m200` in MetaPhlAn2.



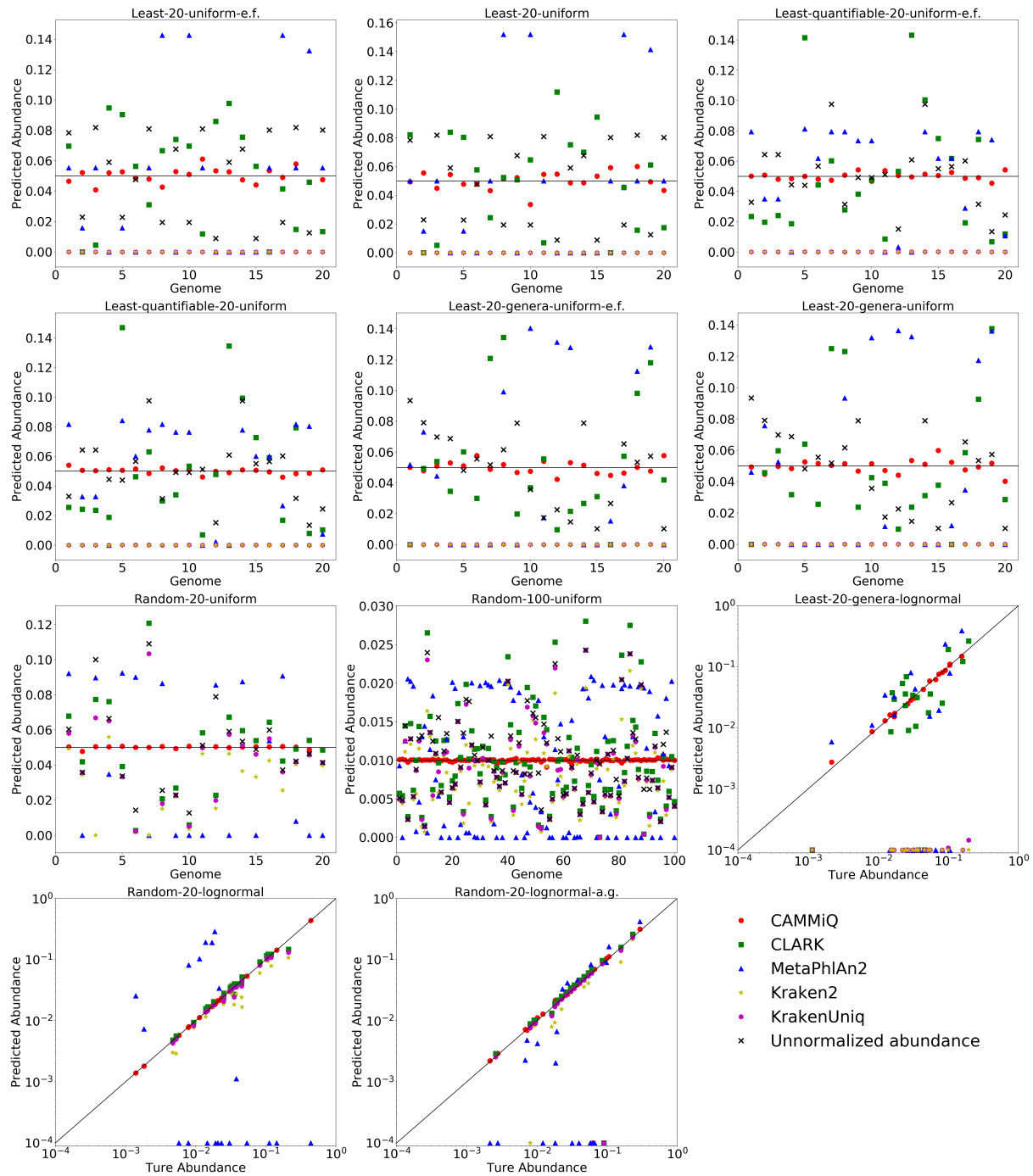


Figure 3: Comparison of CAMMiQ's relative abundance estimates to that of the existing methods for all 11 queries involving the species level dataset. The horizontal axis in the top 8 plots represent the species (i.e., genomes) in an arbitrary order and the vertical axis represent the relative abundance values. The bottom 3 plots are for the queries with log-normal distributions, where the horizontal axis represent the true abundance values while the vertical axis represent the estimated abundance values, both in log-scale.

in Table 4. As can be seen, CAMMiQ managed to identify and accurately quantify all strains in the queries HumanGut-random-100-1 and HumanGut-random-100-2, and  $> 98\%$  strains in the query HumanGut-all, with almost no false positives. Here we only report the strain level calls made by MetaPhlan2 (it made additional calls at the species level or higher), based on our best attempt to match its calls to the strain IDs available to us. It may be possible to slightly increase the true positive values (the first value) reported here.

### 3.5 Performance on Computational Resources

We compared the running time and memory use of CAMMiQ, Kraken2, KrakenUniq, CLARK and MetaPhlan2 in responding to the queries; see Table 5. We do not report the time for building the index or loading it into memory, since this is performed only once - all tools roughly need a couple hours to construct the index on the species level dataset. When it comes to querying time, CAMMiQ is outperformed only by Kraken2, typically by a factor of 3. This is primarily due to the fact that Kraken2’s index is much smaller than that of CAMMiQ since it consists of a small subset of unique  $k$ -mers. This compact index structure substantially impacts its performance, which is improved by KrakenUniq, through its consideration of all unique  $k$ -mers. The resulting index size of KrakenUniq is comparable to that of CAMMiQ (CAMMiQ’s is slightly higher as it also includes doubly-unique substrings), however its running time is almost twice as much, even though it does not solve an ILP. CAMMiQ owes its superior identification and quantification performance to its ILP formulation, however it is not time-wise penalized by it. The best memory footprint is achieved by MetaPhlan2 through the use of its own index, however both its run time and its identification/quantification performance is worse than the others.

### 3.6 Identification Performance - Subspecies Level Dataset with Real scRNAseq Queries

As mentioned earlier, for the “real” metatranscriptomic dataset, we employed query types  $\mathcal{A}_1$  and  $\mathcal{A}_2$  instead of query type  $\mathcal{A}_3$ . We remind the reader that  $\mathcal{A}_1$  simply requires the computation of all genomes in the index dataset for which there is at least one unique substring observed in the query read set, and  $\mathcal{A}_2$  requires the computation of the smallest set of genomes in the index dataset to include all unique or doubly-unique substrings observed in the query read set. Here each query corresponds to the Filtered set of non-human scRNA-seq reads from the corresponding immune cells as described in [58]. More specifically we pre-filtered all query scRNA-seq reads which (i) possibly originate from the human genome, (ii) have low sequence quality and “complexity”, or (iii) map to 16S or 23S ribosomal RNAs on the two *Salmonella* genomes (to avoid incorrect assignment of reads due to “barcode hopping”). When analyzing the resulting genomes, we used annotations provided by [58] and categorized the cells into 5 groups: infected cells that were confirmed to contain intracellular *Salmonella* (either the STM-LT2 or STM-D23580 strain); bystander cells that were exposed to either the STM-LT2 or STM-D23580 strains but confirmed to not contain intracellular *Salmonella*; and cells that were mock-infected and sequenced as controls. For each query, we compared the number of reads CAMMiQ assigned uniquely to STM-LT2 or STM-D23580 genomes to those aligned (and assigned) by GATK PathSeq [10] as well as BLASTN [9].

As expected, CAMMiQ reconfirmed that the abundance (measured by unique read counts) of *Salmonella* was substantially higher in the infected cells compared to the mock-infected controls; and importantly, the unique STM-LT2 (and STM-D23580) reads were differentially abundant between the cells exposed to or infected with STM-LT2 strain (and STM-D23580 strain respectively; see Figure 4A). Interestingly, CAMMiQ identified a small number of unique STM-LT2 reads in the cells exposed to or infected with STM-D23580 strain and vice versa; these were verified with BLASTN, indicating possible sequencing errors or incorrect assignment of specific reads to cells. Compared with the alignment based approach GATK PathSeq with the same index dataset, CAMMiQ was shown to be more sensitive: on average, it identified (roughly) an order of magnitude higher number of unique STM-LT2 or STM-D23580 reads per corresponding cell. This demonstrates CAMMiQ’s potential ability to better identify intracellular organisms at subspecies or strain level; see Figure 4B. <sup>6</sup> When it comes to running time, CAMMiQ offers several orders of magnitude comparative

<sup>6</sup>As can be seen in Figure 4B, even though CAMMiQ’s  $\mathcal{A}_2$  type queries generally improve over its  $\mathcal{A}_1$  type queries with respect to the number of cells correctly identified with STM-LT2 strain or STM-D23580 strains, it seems to also introduce some potential false positive calls (e.g. the last panel in the figure corresponding to the controls). This could be due to additional

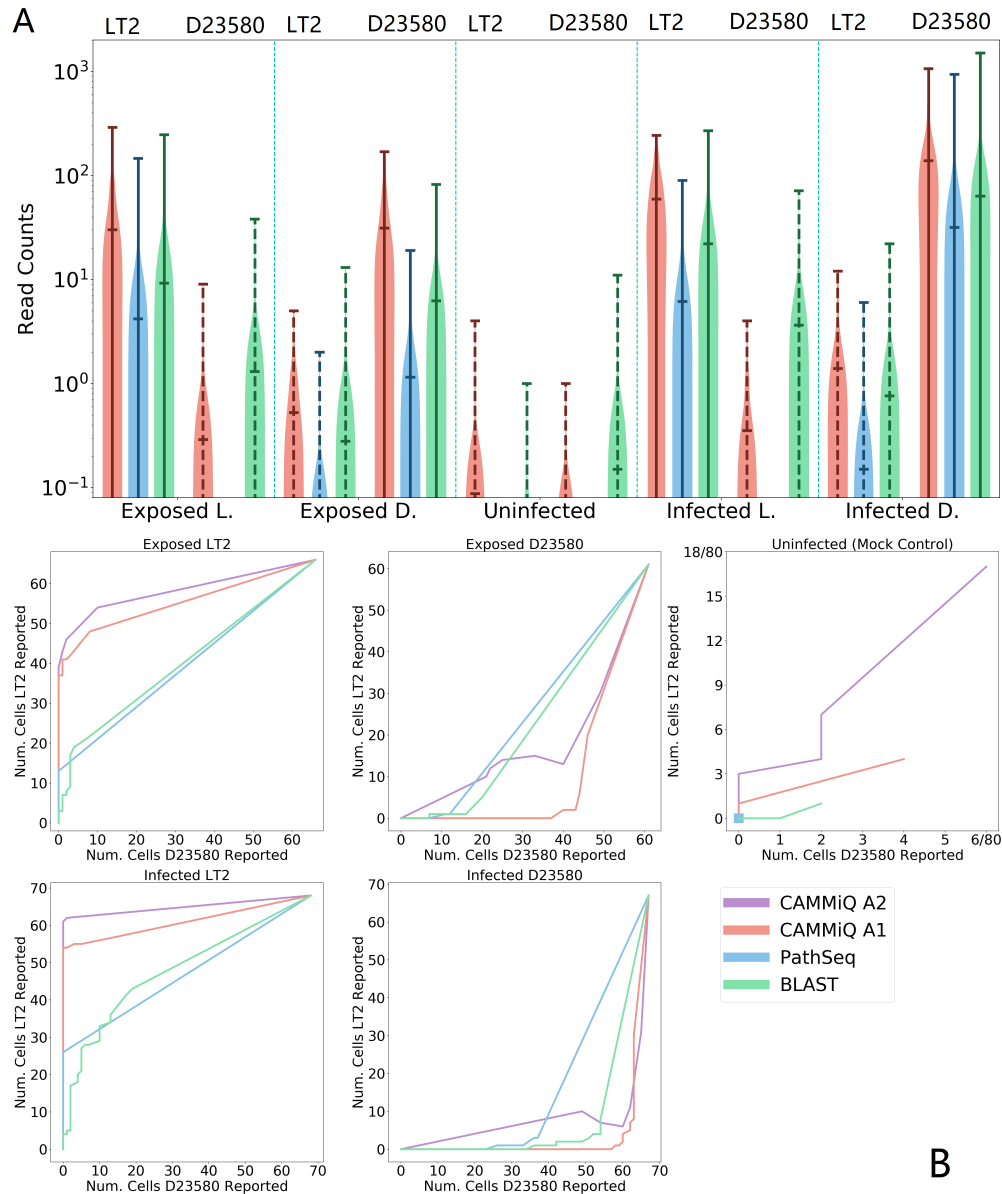


Figure 4: A: Comparison of the distributions of filtered-scrRNAseq reads uniquely assigned to the STM-LT2 (the left side in each panel) or STM-D23580 strain (the right side in each panel) by CAMMiQ (red), GATK PathSeq (blue) and BLASTN (green) between the 5 groups of cells that were (i) exposed to STM-LT2; (ii) exposed to STM-D23580; (iii) mock-infected as controls; (iv) infected with STM-LT2; and (v) infected with STM-D23580. Solid lines indicate the strain which the cells have been exposed to or infected with should have higher read count values. B: The number of cells with more than  $t$  reads uniquely assigned to the STM-LT2 strain (y-axis in each panel) or STM-D23580 strain (x-axis in each panel) by CAMMiQ with query type  $\mathcal{A}_1$  (red), GATK PathSeq (blue) and BLASTN (green), for varying values of  $t$ . We also included the results of CAMMiQ with query type  $\mathcal{A}_2$  (purple) which uses doubly-unique substrings in addition to unique substrings and thus is not represented in part (A) of the figure. When the threshold  $t$  is very high, neither strain could be detected by any method in any of the cells - this corresponds to position (0,0) in each plot. When the threshold is very low (as low as 0), then both strains can be detected in all cells by all tools - this corresponds to position  $(c,c)$  where  $c$  corresponds to the total number of cells in a given plot. The plot for an ideal tool should connect positions (0,0) with  $(c,c)$  following the axis that represents the particular strain with which the cells are exposed to or infected with, as closely as possible (the mock control is the exception). The plots for CAMMiQ are closer to this ideal than PathSeq or BLASTN.

Query Set	CAMMiQ		MetaPhlan2	
	Num. Identified Strains	L1 Err.	Num. Identified Strains	L1 Err.
HumanGut-least-25	24/26	0.1130	$\geq 18/19$	0.2910
HumanGut-random-100-1	100/101	0.0209	$\geq 74/83$	0.3260
HumanGut-random-100-2	100/102	0.0256	$\geq 67/72$	0.5066
HumanGut-all	404/407	0.0517	$\geq 279/305$	0.4439

Table 4: CAMMiQ’s quantification performance at strain level, compared to MetaPhlan2. Number of identified strains: the number of true positive strains/the total number of strains identified. L1 error: the L1 distance between the true relative abundance values and the predicted abundance values, across all strains in the query.

advantage on the filtered-scRNAseq queries. Specifically, CAMMiQ only took a total of 65.3s for computing  $\mathcal{A}_1$  type queries and an additional 2.5s for computing  $\mathcal{A}_2$  type queries on the entire query set, outperforming GATK PathSeq, which required 29628.1s.

## 4 Discussion

We have introduced CAMMiQ, a new computational approach to solve a computational problem that has not been exactly addressed by any available method: given a set  $\mathcal{S}$  of distinct genomic sequences (of any taxonomic rank), build a data structure so as to identify and quantify genomes in a any query, composed of a mixture of reads from a subset of genomes from  $\mathcal{S}$ . CAMMiQ is particularly designed to handle genomes that lack unique features; for that, it reduces the identification and quantification problems to a combinatorial optimization problem that assigns substrings with limited ambiguity (i.e. doubly unique substrings) to genomes so that (in its most general  $\mathcal{A}_3$  type query), each genome is “uniformly covered”. Because each such substring has limited ambiguity, the resulting combinatorial optimization problem can be very efficiently solved through existing integer program solvers such as IBM CPLEX (and thus are branch-and-bound based, whose run-time performance is determined by the “fan-out” at each decision point). Provided that the doubly-unique substrings of a given genome are not all shared with one other genome, the use of doubly-unique substrings increases CAMMiQ’s ability to identify and quantify this genome within a query. As mentioned earlier, in case the dataset to be indexed involves several genomes with high level of similarity, CAMMiQ’s data structure and its combinatorial optimization formulation can easily be generalized to include triply or quadruply unique substrings, without much computational overhead.

## References

- [1] Huttenhower, C. *et al.* Structure, function and diversity of the healthy human microbiome. *nature* **486**, 207 (2012).
- [2] Bullman, S. *et al.* Analysis of fusobacterium persistence and antibiotic response in colorectal cancer. *Science* **358**, 1443–1448 (2017).
- [3] Castellarin, M. *et al.* Fusobacterium nucleatum infection is prevalent in human colorectal carcinoma. *Genome Research* **22**, 299–306 (2012).
- [4] Gur, C. *et al.* Binding of the fap2 protein of fusobacterium nucleatum to human inhibitory receptor tigit protects tumors from immune cell attack. *Immunity* **42**, 344–355 (2015).
- [5] Gur, C. *et al.* Fusobacterium nucleatum supresses anti-tumor immunity by activating ceacam1. *Oncoimmunology* **8**, e1581531 (2019).

reads utilized by  $\mathcal{A}_2$  queries impacted by read errors or misassigned reads.

Performance Measure	Query Set	CAMMiQ	Kraken2	KrakenUniq	CLARK	Metaphlan2
A. Read Assignment Time	Least-20-uniform	151.7s	<b>52.2s</b>	286.5s	402.5s	608.8s
	Least-quantifiable-20-uniform	180.5s	<b>50.8s</b>	284.1s	416.0s	569.6s
	Least-20-genera-uniform	143.5s	<b>53.0s</b>	273.4s	285.2s	438.9s
	Least-20-genera-lognormal	153.6s	<b>61.8s</b>	268.8s	289.4s	442.7s
	Random-20-uniform	198.4s	<b>48.1s</b>	284.0s	315.1s	516.7s
	Random-20-lognormal	167.8s	<b>48.8s</b>	294.5s	367.6s	569.3s
	Random-20-lognormal-a.g.	49.8s	<b>14.7s</b>	78.9s	75.2s	188.9s
	Random-100-uniform	906.0s	<b>215.4s</b>	1261.5s	1536.2s	2401.4s
B. Memory	Least-20-uniform	94.6G	11.4G	67.9G	81.3G	<b>1.3G</b>
	Least-quantifiable-20-uniform	93.3G	11.4G	77.1G	81.3G	<b>1.3G</b>
	Least-20-genera-uniform	92.6G	11.4G	81.4G	81.0G	<b>1.3G</b>
	Least-20-genera-lognormal	93.3G	11.4G	77.8G	81.0G	<b>1.3G</b>
	Random-20-uniform	95.4G	11.4G	92.7G	81.3G	<b>1.3G</b>
	Random-20-lognormal	95.7G	11.4G	89.3G	81.3G	<b>1.3G</b>
	Random-20-lognormal-a.g.	94.3G	11.3G	51.9G	80.3G	<b>1.3G</b>
	Random-100-uniform	111.9G	11.9G	143.3G	85.7G	<b>1.4G</b>
C. Index Size	-	8.94G	11.20G	155.22G	71.85G	<b>1.11G</b>

Table 5: Comparison of the computational resources required by CAMMiQ, Kraken2, KrakenUniq, CLARK and MetaPhlan2, when run on a single thread.

- [6] Kostic, A. D. *et al.* Genomic analysis identifies association of fusobacterium with colorectal carcinoma. *Genome Research* **22**, 292–298 (2012).
- [7] Yu, T. *et al.* Fusobacterium nucleatum promotes chemoresistance to colorectal cancer by modulating autophagy. *Cell* **170**, 548–563 (2017).
- [8] Huson, D. H., Auch, A. F., Qi, J. & Schuster, S. C. Megan analysis of metagenomic data. *Genome Research* **17**, 377–386 (2007).
- [9] Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool. *Journal of Molecular Biology* **215**, 403–410 (1990).
- [10] Walker, M. A. *et al.* Gatk pathseq: a customizable computational tool for the discovery and identification of microbial sequences in libraries from eukaryotic hosts. *Bioinformatics* **34**, 4287–4289 (2018).
- [11] Poore, G. D. *et al.* Microbiome analyses of blood and tissues suggest cancer diagnostic approach. *Nature* **579**, 567–574 (2020).
- [12] Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biology* **15**, R46 (2014).
- [13] Elworth, R. *et al.* To petabytes and beyond: recent advances in probabilistic and signal processing algorithms and their application to metagenomics. *Nucleic Acids Research* (2020).
- [14] Robinson, W., Schischlik, F., Gertz, E. M., Schäffer, A. A. & Ruppig, E. Identifying the landscape of intratumoral microbes via a single cell transcriptomic analysis. *bioRxiv* (2020).
- [15] Liu, B., Gibbons, T., Ghodsi, M., Treangen, T. & Pop, M. Accurate and fast estimation of taxonomic profiles from metagenomic shotgun sequences. *Genome Biology* **12** (Suppl 2), S4 (2011).
- [16] Segata, N. *et al.* Metagenomic microbial community profiling using unique clade-specific marker genes. *Nature Methods* **9**, 811 (2012).



- [17] Lopresti, D. & Tomkins, A. Block edit models for approximate string matching. *Theoretical Computer Science* **181**, 159–179 (1997).
- [18] Cormode, G., Paterson, M., Sahinalp, S. C. & Vishkin, U. Communication complexity of document exchange. In *Proceedings of the eleventh annual ACM-SIAM symposium on Discrete algorithms*, 197–206 (Society for Industrial and Applied Mathematics, 2000).
- [19] Li, M., Chen, X., Li, X., Ma, B. & Vitányi, P. M. The similarity metric. *IEEE Transactions on Information Theory* **50**, 3250–3264 (2004).
- [20] Vinga, S. & Almeida, J. Alignment-free sequence comparison—a review. *Bioinformatics* **19**, 513–523 (2003).
- [21] Leimeister, C.-A. & Morgenstern, B. Kmacs: the k-mismatch average common substring approach to alignment-free sequence comparison. *Bioinformatics* **30**, 2000–2008 (2014).
- [22] Ames, S. K. *et al.* Scalable metagenomic taxonomy classification using a reference genome database. *Bioinformatics* **29**, 2253–2260 (2013).
- [23] Břinda, K., Sykulski, M. & Kucherov, G. Spaced seeds improve k-mer-based metagenomic classification. *Bioinformatics* **31**, 3584–3592 (2015).
- [24] Kawulok, J. & Deorowicz, S. Cometa: classification of metagenomes using k-mers. *PloS ONE* **10**, e0121453 (2015).
- [25] Tu, Q., He, Z. & Zhou, J. Strain/species identification in metagenomes using genome-specific markers. *Nucleic Acids Research* **42**, e67–e67 (2014).
- [26] Ounit, R., Wanamaker, S., Close, T. J. & Lonardi, S. Clark: fast and accurate classification of metagenomic and genomic sequences using discriminative k-mers. *BMC Genomics* **16**, 236 (2015).
- [27] Luo, Y., Zeng, J., Berger, B. & Peng, J. Low-density locality-sensitive hashing boosts metagenomic binning. In *International Conference on Research in Computational Molecular Biology*, 255 (Springer, 2016).
- [28] Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using minhash. *Genome Biology* **17**, 132 (2016).
- [29] McHardy, A. C., Martín, H. G., Tsirigos, A., Hugenholtz, P. & Rigoutsos, I. Accurate phylogenetic classification of variable-length dna fragments. *Nature Methods* **4**, 63 (2007).
- [30] Rosen, G., Garbarine, E., Caseiro, D., Polikar, R. & Sokhansanj, B. Metagenome fragment classification using n-mer frequency profiles. *Advances in Bioinformatics* **2008** (2008).
- [31] Brady, A. & Salzberg, S. L. Phymm and phymmbl: metagenomic phylogenetic classification with interpolated markov models. *Nature Methods* **6**, 673 (2009).
- [32] Rosen, G. L., Reichenberger, E. R. & Rosenfeld, A. M. Nbc: the naive bayes classification tool webserver for taxonomic classification of metagenomic reads. *Bioinformatics* **27**, 127–129 (2010).
- [33] Vervier, K., Mahé, P., Tournoud, M., Veyrieras, J.-B. & Vert, J.-P. Large-scale machine learning for metagenomics sequence classification. *Bioinformatics* **32**, 1023–1032 (2015).
- [34] Anyansi, C., Straub, T. J., Manson, A. L., Earl, A. M. & Abeel, T. Computational methods for strain-level microbial detection in colony and metagenome sequencing data. *Frontiers in Microbiology* **11**, 1925 (2020).
- [35] Marshall, J. A. Mixed infections of intestinal viruses and bacteria in humans. In *Polymicrobial Diseases* (ASM Press, 2002).

- [36] Cohen, T. *et al.* Mixed-strain mycobacterium tuberculosis infections and the implications for tuberculosis treatment and control. *Clinical microbiology reviews* **25**, 708–719 (2012).
- [37] Balmer, O. & Tanner, M. Prevalence and implications of multiple-strain infections. *The Lancet infectious diseases* **11**, 868–878 (2011).
- [38] Secher, T., Brehin, C. & Oswald, E. Early settlers: which e. coli strains do you not want at birth? *American Journal of Physiology: Gastrointestinal and Liver Physiology* **311**, G123–G129 (2016).
- [39] Gerner-Smidt, P. *et al.* Whole genome sequencing: Bridging one-health surveillance of foodborne diseases. *Frontiers in Public Health* **7**, 172 (2019).
- [40] Lin, Y.-Y. *et al.* Cliiq: Accurate comparative detection and quantification of expressed isoforms in a population. In *International Workshop on Algorithms in Bioinformatics*, 178–189 (Springer, 2012).
- [41] Li, W., Feng, J. & Jiang, T. Isolasso: a lasso regression approach to rna-seq based transcriptome assembly. *Journal of Computational Biology* **18**, 1693–1707 (2011).
- [42] Dao, P. *et al.* Orman: optimal resolution of ambiguous rna-seq multimappings in the presence of novel isoforms. *Bioinformatics* **30**, 644–651 (2014).
- [43] Sobih, A., Tomescu, A. I. & Mäkinen, V. Metaflow: Metagenomic profiling based on whole-genome coverage analysis with min-cost flows. In *International Conference on Research in Computational Molecular Biology*, 111–121 (Springer, 2016).
- [44] Breitwieser, F., Baker, D. & Salzberg, S. L. Krakenuniq: confident and fast metagenomics classification using unique k-mer counts. *Genome Biology* **19**, 1–10 (2018).
- [45] Solomon, B. & Kingsford, C. Fast search of thousands of short-read sequencing experiments. *Nature Biotechnology* **34**, 300 (2016).
- [46] Solomon, B. & Kingsford, C. Improved search of large transcriptomic sequencing databases using split sequence bloom trees. In *International Conference on Research in Computational Molecular Biology*, 257–271 (Springer, 2017).
- [47] Sun, C., Harris, R. S., Chikhi, R. & Medvedev, P. Allsome sequence bloom trees. In *International Conference on Research in Computational Molecular Biology*, 272–286 (Springer, 2017).
- [48] Ferdman, M., Johnson, R. & Patro, R. Mantis: A fast, small, and exact large-scale sequence-search index. In *Research in Computational Molecular Biology*, 271 (Springer, 2018).
- [49] Nazeen, S., Yu, Y. W. & Berger, B. Carnelian uncovers hidden functional patterns across diverse study populations from whole metagenome sequencing reads. *Genome biology* **21**, 1–18 (2020).
- [50] Danko, D. C., Meleshko, D., Bezdán, D., Mason, C. & Hajirasouliha, I. Minerva: an alignment-and reference-free approach to deconvolve linked-reads for metagenomics. *Genome research* **29**, 116–124 (2019).
- [51] Matias, Y., Muthukrishnan, S., Sahinalp, S. C. & Ziv, J. Augmenting suffix trees, with applications. In *European Symposium on Algorithms*, 67–78 (Springer, 1998).
- [52] Kasai, T., Lee, G., Arimura, H., Arikawa, S. & Park, K. Linear-time longest-common-prefix computation in suffix arrays and its applications. In *Annual Symposium on Combinatorial Pattern Matching*, 181–192 (Springer, 2001).
- [53] Kärkkäinen, J. & Sanders, P. Simple linear work suffix array construction. In *International Colloquium on Automata, Languages, and Programming*, 943–955 (Springer, 2003).
- [54] Ilie, L. & Smyth, W. F. Minimum unique substrings and maximum repeats. *Fundamenta Informaticae* **110**, 183–195 (2011).

- [55] Karp, R. M. & Rabin, M. O. Efficient randomized pattern-matching algorithms. *IBM journal of research and development* **31**, 249–260 (1987).
- [56] Vazirani, V. V. *Approximation algorithms* (Springer Science & Business Media, 2013).
- [57] Pruitt, K. D., Tatusova, T. & Maglott, D. R. Ncbi reference sequences (refseq): a curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Research* **35**, D61–D65 (2007).
- [58] Aulicino, A. *et al.* Invasive salmonella exploits divergent immune evasion strategies in infected and bystander dendritic cell subsets. *Nature Communications* **9**, 4883 (2018).
- [59] Wood, D. E., Lu, J. & Langmead, B. Improved metagenomic analysis with kraken 2. *Genome Biology* **20**, 257 (2019).
- [60] Truong, D. T. *et al.* Metaphlan2 for enhanced metagenomic taxonomic profiling. *Nature Methods* **12**, 902 (2015).
- [61] Forster, S. C. *et al.* A human gut bacterial genome and culture collection for improved metagenomic analyses. *Nature Biotechnology* **37**, 186 (2019).
- [62] Weissman, T., Ordentlich, E., Seroussi, G., Verdu, S. & Weinberger, M. J. Inequalities for the L1 deviation of the empirical distribution. *Hewlett-Packard Labs, Tech. Rep* (2003).
- [63] Leinonen, R., Sugawara, H., Shumway, M. *et al.* The sequence read archive. *Nucleic Acids Research* **39**, D19–D21 (2011).
- [64] Morgulis, A. & Agarwala, R. Srprism (single read paired read indel substitution minimizer): an efficient aligner for assemblies with explicit guarantees. *GigaScience* **9**, giaa023 (2020).
- [65] Altschul, S. F. *et al.* Gapped blast and psi-blast: a new generation of protein database search programs. *Nucleic Acids Research* **25**, 3389–3402 (1997).

## 5 Supplementary Methods

### 5.1 Unique substrings from $LCP_u$ or $LCP_d$

The pseudocode for the algorithm to compute arrays **SU** and **SD** according to equations (3) and (4) respectively, given **SA**,  $LCP_u$  and  $LCP_d$  is given below. As can be seen, the algorithm runs in  $O(M)$  time.

---

**Algorithm 1** ShortestUniqueFromLCP(**SA**,  $LCP_u$ ,  $LCP_d$ ,  $L_{\max}$ )

---

```

1: for  $i$  from 1 to  $M$  do //Initialize SU
2:   SU[ $i$ ]  $\leftarrow$  0
3:   SD[ $i$ ]  $\leftarrow$  0
4: end for
5: for  $i$  from 1 to  $M$  do //Update SU according to (3) and SD according to (4)
6:   if  $LCP_u[i] < L_{\max}$  then
7:     SU[SA[ $i$ ] +  $LCP_u[i]$ ]  $\leftarrow$   $\max\{\mathbf{SU}[\mathbf{SA}[i] + LCP_u[i]], \mathbf{SA}[i]\}$ 
8:   end if
9:   if  $LCP_d[i] < L_{\max}$  then
10:    SD[SA[ $i$ ] +  $LCP_d[i]$ ]  $\leftarrow$   $\max\{\mathbf{SD}[\mathbf{SA}[i] + LCP_d[i]], \mathbf{SA}[i]\}$ 
11:   end if
12: end for
13: return SU, SD

```

---

### 5.2 Computing $LCP_u$ and $LCP_d$

In this section we show that **SU** and **SD** can be correctly constructed in  $O(M)$  time. We start by showing that the definition of  $LCP_u$  and  $LCP_d$  in Equation (1) and Equation (2) can respectively lead to the shortest substrings occurring in at most one genome or two genomes. Then we give CAMMIQ's detailed implementation of Equation (5) and Equation (6) to compute the  $LCP_u$  and  $LCP_d$  arrays. Finally we give a running time analysis of this implementation.

First consider the content of **SU** at the end of procedure ShortestUniqueFromLCP. To see the substring  $s[l : r]$  corresponds to the  $r$ -th entry  $\mathbf{SU}[r] = l$  (where  $l \neq 0$ ) in **SU** is unique, meaning it only occurs in genome with ID  $\mathbf{GSA}[\mathbf{SA}^{-1}[l]]$ , assume that there is another genome  $s_{i'}$  having the same substring  $s[l', r'] = s[l : r]$  - this leads to a contradiction, since it implies that  $\text{lcp}(\mathbf{suf}[l], \mathbf{suf}[l']) \geq r - l + 1$ . However, due to the update rule of **SU** and the definition of  $LCP_u$ ,  $\text{lcp}(\mathbf{suf}[l], \mathbf{suf}[l']) \leq r - l$  for any  $1 \leq l' \leq M$  satisfying  $\mathbf{suf}[l]$  and  $\mathbf{suf}[l']$  start on different genomes, namely  $\mathbf{GSA}[\mathbf{SA}^{-1}[l]] \neq \mathbf{GSA}[\mathbf{SA}^{-1}[l']]$ , which is a contradiction. Now, to see  $s[l : r]$  is a shortest unique substring, i.e. no substring of  $s[l : r]$  is unique to genome  $\mathbf{GSA}[\mathbf{SA}^{-1}[l]]$ , we show that any  $s[l : r' < r]$  and  $s[l' > l, r]$  occurs in one other genome  $s_{i'}$ . The former case is due to the definition of  $LCP_u$  - there exists  $\mathbf{suf}[l']$  on genome  $i' \neq \mathbf{GSA}[\mathbf{SA}^{-1}[l]]$  such that  $\text{lcp}(\mathbf{suf}[l], \mathbf{suf}[l']) \geq r - l$ , implying a substring  $s[l' : l' + (r - l) - 1]$  identical to  $s[l, r - 1]$ ; the later case is due to the update rule of **SU** - if  $s[l' > l, r]$  is also unique to genome  $\mathbf{GSA}[\mathbf{SA}^{-1}[l]]$ , then  $\mathbf{SU}[r]$  must be set to  $l'$  instead of  $l$ . Therefore,  $s[l : r]$  is a shortest unique substring (to genome with ID  $\mathbf{GSA}[\mathbf{SA}^{-1}[l]]$ ); on the other hand, if  $s[l : r]$  is a shortest unique substring, then  $\mathbf{SU}[r]$  will maintain  $l$  after **SU** is completely updated.

Now consider the content of **SD** at the end of procedure ShortestUniqueFromLCP. We follow the above proof to show  $s[l : r]$  is a shortest doubly-unique substring (to genome ID  $\mathbf{GSA}[\mathbf{SA}^{-1}[l]]$  and possibly another genome  $i'$ ). To see the substring  $s[l : r]$  corresponds to the  $r$ -th entry  $\mathbf{SD}[r] = l$  (where  $l \neq 0$ ) in **SD** occurs in at most two genomes, with ID  $\mathbf{GSA}[\mathbf{SA}^{-1}[l]]$  (and possibly  $i'$ , any genome that  $\mathbf{suf}[l']$  belongs to, giving the largest  $\text{lcp}(\mathbf{suf}[l], \mathbf{suf}[l'])$ ), we can assume there exists a third genome  $s_{i''}$  having the same substring  $s[l'', r''] = s[l, r] = s[l', r']$  and similarly obtain a contradiction. Note that according to (2), it is possible to have  $LCP_d[\mathbf{SA}^{-1}[l]] \geq LCP_u[\mathbf{SA}^{-1}[l]]$  and in this case  $s[l : r]$  is a unique substring which only occurs in genome  $\mathbf{GSA}[\mathbf{SA}^{-1}[l]]$ . If  $LCP_d[\mathbf{SA}^{-1}[l]] < LCP_u[\mathbf{SA}^{-1}[l]]$  on the other hand, then  $s[l : r]$  must occur in exactly two genomes, since **SD** is updated according to  $LCP_d$  and we can find another suffix of  $s$  whose length- $(r - l + 1)$  ( $r - l + 1 \leq LCP_d[\mathbf{SA}^{-1}[l]]$ ) prefix is identical to  $s[l : r]$ . In addition, to see  $s[l : r]$  is a shortest doubly-unique substring, meaning no substring of  $s[l : r]$  occurs only in genome  $\mathbf{GSA}[\mathbf{SA}^{-1}[l]]$  and  $i'$ , we can

similarly show that any  $s[l : r' < r]$  and  $s[l' > l, r]$  can be found in a third genome  $s_{i''}$ , regardless whether  $\text{LCP}_d[\text{SA}^{-1}[l]] = \text{LCP}_u[\text{SA}^{-1}[l]]$  or  $\text{LCP}_d[\text{SA}^{-1}[l]] < \text{LCP}_u[\text{SA}^{-1}[l]]$ .

As a result of the above observations we can now formally state the following lemma.

**Lemma 2.** (i) After updating  $SU$  according to (3),  $SU[r] = l \neq 0$  implies that  $s[l : r]$  is a shortest unique substring to genome  $GSA[\text{SA}^{-1}[l]]$ ;

(ii) After updating  $SD$  according to (4)  $SD[r] = l \neq 0$  implies that  $s[l : r]$  is a shortest doubly-unique substring to genome  $GSA[\text{SA}^{-1}[l]]$  and  $i' = GSA[\text{SA}^{-1}[l']]$  where  $\text{su}f[l']$  gives the largest  $\text{lcp}(\text{su}f[l], \text{su}f[l'])$ .

Furthermore, it's also clear that all unique substrings  $s[l : r]$  are stored in  $SU$  and all doubly-unique substrings  $s[l : r]$  are stored in  $SD$ , as we have considered the suffix of  $s$  starting with every possible  $l$ .

Both  $\text{LCP}_u$  and  $\text{LCP}_d$  can be modified to incorporate the minimum length constraint  $L_{\min}$  on unique/doubly-unique substrings. By setting  $\text{LCP}_u[i]$  to  $\max\{L_{\min} - 1, \text{LCP}_u[i]\}$  for each entry  $1 \leq i \leq M$ , the corresponding substrings maintained in  $SU$  should also be unique, and with minimum length  $L_{\min}$ . One should be careful however when dealing with  $\text{LCP}_d$ : if  $\text{LCP}_u[i] \leq L_{\min} - 1$ , then the corresponding substring  $s[\text{SA}[i], \text{SA}[i] + \text{LCP}_u[i]]$  occurs in only one genome. Therefore we set  $\text{LCP}_d[i]$  to  $\infty$  (meaning it's not considered) if  $\text{LCP}_u[i] \leq L_{\min} - 1$  or  $\text{LCP}_d[i] \geq \text{LCP}_u[i]$ ; and to  $\max\{L_{\min} - 1, \text{LCP}_d[i]\}$  otherwise (this can be done by first set each  $\text{LCP}_u[i]$  to  $\max\{L_{\min} - 1, \text{LCP}_u[i]\}$  and  $\text{LCP}_d[i]$  to  $\max\{L_{\min} - 1, \text{LCP}_d[i]\}$ , and then set each  $\text{LCP}_d[i]$  to  $\infty$  if  $\text{LCP}_d[i] \geq \text{LCP}_u[i]$ ), which ensures the corresponding substrings maintained in  $SD$  are doubly-unique, and with minimum length  $L_{\min}$ .

With the correctness of (3) and (4) in mind, our next concern is how to actually compute  $\text{LCP}_u$  and  $\text{LCP}_d$  based on their definitions. In the following we show that (5) and (6) correctly implement (1) and (2), without considering the borderline cases (i.e., for  $i = 1$  or  $i = M$ ; to handle these cases we can set  $GSA[0] = GSA[M + 1] = 0$  and ignore  $i^{2-}$  when  $i = 1$  and  $i^{2+}$  when  $i = M$ ).

**Lemma 3.** For any  $1 \leq i \leq M$ ,

$$(i) \text{LCP}_u[i] = \max \begin{cases} \min_{i^- < x \leq i} \text{LCP}[x], \text{ where } i^- = \max\{1 \leq i' < i\}, \text{ s.t. } d_{GSA}(i', i) \geq 2 \\ \min_{i < x \leq i^+} \text{LCP}[x], \text{ where } i^+ = \min\{i < i' \leq M\}, \text{ s.t. } d_{GSA}(i, i') \geq 2 \end{cases}$$

$$(ii) \text{LCP}_d[i] = \min \begin{cases} \max \begin{cases} \min_{i^- < x \leq i} \text{LCP}[x], \text{ where } i^- = \max\{1 \leq i' < i\}, \text{ s.t. } d_{GSA}(i', i) \geq 2 \\ \min_{i < x \leq i^{2+}} \text{LCP}[x], \text{ where } i^{2+} = \min\{i < i' \leq M\}, \text{ s.t. } d_{GSA}(i, i') \geq 3 \end{cases} \\ \max \begin{cases} \min_{i^{2-} < x \leq i} \text{LCP}[x], \text{ where } i^{2-} = \max\{1 \leq i' < i\}, \text{ s.t. } d_{GSA}(i', i) \geq 3 \\ \min_{i < x \leq i^+} \text{LCP}[x], \text{ where } i^+ = \min\{i < i' \leq M\}, \text{ s.t. } d_{GSA}(i, i') \geq 2 \end{cases} \end{cases}$$

where  $d_{GSA}(i_1, i_2) = |\{GSA[i_1], \dots, GSA[i_2]\}|$ .

*Proof.* Let  $\text{lcp}(i, j)$  denote  $\text{lcp}(\text{su}f[i], \text{su}f[j])$  for short. We utilize the properties of  $\text{SA}$  and  $\text{LCP}$  array: (a) the longest common prefix of two suffices  $\text{su}f[i]$  and  $\text{su}f[j]$  (assume  $\text{su}f[i]$  is lexicographically smaller than  $\text{su}f[j]$ ) is  $\text{lcp}(i, j) = \min\{\text{LCP}[x] \mid x \in [\text{SA}^{-1}[i] + 1, \text{SA}^{-1}[j]]\}$ ; also we have (b)  $\text{lcp}(i, j) \geq \text{lcp}(\text{SA}[\text{SA}^{-1}[i] - 1], j)$  and  $\text{lcp}(i, j) \geq \text{lcp}(i, \text{SA}[\text{SA}^{-1}[j] + 1])$ . (i) follows immediately from these properties:

$$\text{LCP}_u[i] = \max\{\text{lcp}(\text{SA}[i], \text{SA}[i^+]), \text{lcp}(\text{SA}[i^-], \text{SA}[i])\}.$$

To see (ii), we consider three cases:

- If  $\text{lcp}(\text{SA}[i], \text{SA}[i^+]) = \text{lcp}(\text{SA}[i^-], \text{SA}[i])$ , then  $\text{LCP}_d[i] = \text{lcp}(\text{SA}[i], \text{SA}[i^+]) = \text{lcp}(\text{SA}[i^-], \text{SA}[i])$ , due to (b). (ii) holds in this case by applying (a).
- If  $\text{lcp}(\text{SA}[i], \text{SA}[i^+]) < \text{lcp}(\text{SA}[i^-], \text{SA}[i])$ , then  $\text{LCP}_d[i] = \max\{\text{lcp}(\text{SA}[i], \text{SA}[i^+]), \text{lcp}(\text{SA}[i^{2-}], \text{SA}[i])\}$  due to (b). Also  $\text{LCP}_d[i] \leq \text{lcp}(\text{SA}[i^-], \text{SA}[i]) = \max\{\text{lcp}(\text{SA}[i], \text{SA}[i^{2+}]), \text{lcp}(\text{SA}[i^-], \text{SA}[i])\}$  since  $\text{lcp}(\text{SA}[i], \text{SA}[i^{2+}]) \leq \text{lcp}(\text{SA}[i], \text{SA}[i^+]) < \text{lcp}(\text{SA}[i^-], \text{SA}[i])$ . (ii) therefore holds by applying (a).



- If  $\text{lcp}(\text{SA}[i], \text{SA}[i^+]) > \text{lcp}(\text{SA}[i^-], \text{SA}[i])$ , then we have similarly  $\text{LCP}_d[i] = \max\{\text{lcp}(\text{SA}[i], \text{SA}[i^{2+}]), \text{lcp}(\text{SA}[i^-], \text{SA}[i])\}$  due to (b) and  $\text{LCP}_d[i] \leq \text{lcp}(\text{SA}[i], \text{SA}[i^+]) = \max\{\text{lcp}(\text{SA}[i], \text{SA}[i^+]), \text{lcp}(\text{SA}[i^{2-}], \text{SA}[i])\}$  since  $\text{lcp}(\text{SA}[i^{2-}], \text{SA}[i]) \leq \text{lcp}(\text{SA}[i^-], \text{SA}[i]) < \text{lcp}(\text{SA}[i], \text{SA}[i^+])$ . Again we can see (ii) by applying (a), which complete the proof. □

Next we give the pseudocode for the algorithm to update  $\text{LCP}_u$  and  $\text{LCP}_d$  based on (5) and (6) respectively, and show that they actually run in linear time.

---

**Algorithm 2**  $\text{LCP}_u\text{FromLCP}(\text{GSA}, \text{LCP}, L_{\min})$

---

<pre> 1: <math>\text{LCP}_u \leftarrow</math> array of length <math>M</math> 2: <math>i \leftarrow 1</math>, <math>\text{minlcp} \leftarrow M</math>, <math>\text{next} \leftarrow 0</math> 3: <b>while</b> <math>i &lt; M</math> <b>do</b> 4:   <b>while</b> <math>i &lt; M</math> and <math>\text{GSA}[i + \text{next}] = \text{GSA}[i +</math>      <math>\text{next} + 1]</math> <b>do</b> 5:     <math>\text{next} \leftarrow \text{next} + 1</math> 6:   <b>end while</b> 7:   <b>for</b> <math>j</math> <b>from</b> <math>\text{next}</math> <b>to</b> <math>0</math> <b>do</b> 8:     <math>\text{minlcp} \leftarrow \min\{\text{minlcp}, \text{LCP}[i + j + 1]\}</math> 9:     <math>\text{LCP}_u[i + j] \leftarrow \text{minlcp}</math> 10:  <b>end for</b> 11:  <math>i \leftarrow i + \text{next} + 1</math> 12:  <math>\text{minlcp} \leftarrow M</math>, <math>\text{next} \leftarrow 0</math> 13: <b>end while</b> 14: <math>i \leftarrow M</math>, <math>\text{minlcp} \leftarrow M</math>, <math>\text{next} \leftarrow 0</math> 15: <b>while</b> <math>i &gt; 1</math> <b>do</b> </pre>	<pre> 16:   <b>while</b> <math>i &gt; 1</math> and <math>\text{GSA}[i - \text{next}] = \text{GSA}[i - \text{next} -</math>      <math>1]</math> <b>do</b> 17:     <math>\text{next} \leftarrow \text{next} + 1</math> 18:   <b>end while</b> 19:   <b>for</b> <math>j</math> <b>from</b> <math>\text{next}</math> <b>to</b> <math>0</math> <b>do</b> 20:     <math>\text{minlcp} \leftarrow \min\{\text{minlcp}, \text{LCP}[i - j]\}</math> 21:     <math>\text{LCP}_u[i - j] \leftarrow \max\{\text{LCP}_u[i - j], \text{minlcp}\}</math> 22:   <b>end for</b> 23:   <math>i \leftarrow i - \text{next} - 1</math> 24:   <math>\text{minlcp} \leftarrow M</math>, <math>\text{next} \leftarrow 0</math> 25: <b>end while</b> 26: <b>for</b> <math>i</math> <b>from</b> <math>1</math> <b>to</b> <math>M</math> <b>do</b> 27:   <math>\text{LCP}_u[i] \leftarrow \max\{\text{LCP}_u[i], L_{\min}\}</math> 28: <b>end for</b> 29: <b>return</b> <math>\text{LCP}_u</math> </pre>
---	---

---

**Lemma 4.** Both  $\text{LCP}_u\text{FromLCP}$  and  $\text{LCP}_d\text{FromLCP}$  run in  $O(M)$  time.

*Proof.* Through a simple aggregate analysis, we can see that  $\text{LCP}_u\text{FromLCP}$  visits each entry of  $\text{GSA}$ ,  $\text{LCP}$  and  $\text{LCP}_u$  2 times;  $\text{LCP}_d\text{FromLCP}^*$  visits each entry of  $\text{GSA}$ ,  $\text{LCP}$  3 times and each entry of  $\text{LCP}_d^*$  2 times for either focus =  $+/-$ . □

Combining the above lemmata, we conclude that

**Theorem 5.** Both  $SU$  and  $SD$  can be computed in  $O(M)$  time.

### 5.3 Sampling unique substrings

Recall that  $\mathcal{U}_i = \{u_{i,1}(l_1, r_1), u_{i,nu_i}(l_{nu_i}, r_{nu_i})\}$  defines either the collection of all shortest unique substrings or unique substrings with minimum length  $L_{\min}$  on a given genome  $s_i$  (sorted by  $l_x$ , namely  $l_1 \leq \dots \leq l_{nu_i}$ ). In fact, the list of left indices  $l_1, \dots, l_{nu_i}$  are stored in the corresponding  $r_1, \dots, r_{nu_i}$  entries in  $SU$  array. Due to the minimum length constraint, no substring  $u_{i,x} \in \mathcal{U}_i$  can be a substring of any other  $u_{i,y} \in \mathcal{U}_i$  if they are not identical (in fact, there could be some  $\{u_{i,x}(l_x, r_x) = \{u_{i,y}(l_y, r_y) \in \mathcal{U}_i$  for  $x \neq y$ ). This makes  $l_1 < \dots < l_{nu_i}$  and  $r_1 < \dots < r_{nu_i}$ . The goal of sampling unique substrings from  $\mathcal{U}_i$  is to identify and maintain the smallest number of unique substrings such that they cover the same set of unique  $L$ -mers on  $s_i$  as  $\mathcal{U}_i$  (if  $L_{\max} = L$  then the sampled unique substrings should cover all unique  $L$ -mers).

Here we present the greedy sampling strategy implemented by CAMMIQ to sample unique substrings from  $\mathcal{U}_i$ . Denote by  $\text{begin}_i$  the beginning position of  $s_i$  in  $s$  and by  $\mathcal{U}'_i$  the unique substrings already sampled (initially  $\mathcal{U}'_i$  is empty). Starting with  $\text{begin}_i$ , consider every  $L$ -mer of  $s_i$  from left to right; if it does not include any unique substring, then ignore this  $L$ -mer; otherwise add its rightmost unique substring into  $\mathcal{U}'_i$  and move to the next  $L$ -mer which does not include this substring until reaching the  $L$ -mer that ends at

---

**Algorithm 3**  $LCP_dFromLCP(GSA, LCP, L_{min})$

---

```

1:  $LCP_d \leftarrow$  array of length  $M$ 
2:  $LCP_u \leftarrow LCP_uFromLCP(GSA, LCP, M, L_{min})$ 
3:  $LCP_{d'} \leftarrow LCP_dFromLCP^*(GSA, LCP, +)$ 
4:  $LCP_{d''} \leftarrow LCP_dFromLCP^*(GSA, LCP, -)$ 
5: for  $i$  from 1 to  $M$  do
6:    $LCP_d[i] \leftarrow \max\{\min\{LCP_{d'}[i], LCP_{d''}[i]\}, L_{min}\}$ 
7:   if  $LCP_d[i] \geq LCP_u[i]$  then
8:      $LCP_d[i] \leftarrow \infty$ 
9:   end if
10: end for
11: return  $LCP_d$ 
12: function  $LCP_dFromLCP^*(GSA, LCP, focus=$ 
     $\{+, -\})$ 
13:    $LCP_d^* \leftarrow$  array of length  $M$ 
14:    $i \leftarrow 1, minlcp \leftarrow M, next1 \leftarrow 0$ 
15:   if  $focus = +$  then
16:      $next2 \leftarrow 1$ 
17:   else
18:      $next2 \leftarrow 0$ 
19:   end if
20:   while  $i < M$  do
21:     while  $i < M$  and  $GSA[i + next1] = GSA[i +$ 
     $next1 + 1]$  do
22:        $next1 \leftarrow next1 + 1$ 
23:     end while
24:     while  $focus = +$  and  $i < M$  and  $GSA[i +$ 
     $next1 + next2] = GSA[i + next1 + next2 + 1]$  do
25:        $next2 \leftarrow next2 + 1$ 
26:     end while
27:     for  $j$  from  $next1 + next2$  to 0 do
28:        $minlcp \leftarrow \min\{minlcp, LCP[i + j + 1]\}$ 
29:       if  $j \leq next1$  then
30:          $LCP_d^*[i + j] \leftarrow minlcp$ 
31:       end if
32:     end for
33:      $i \leftarrow i + next1 + 1$ 
34:      $minlcp \leftarrow M, next1 \leftarrow 0$ 
35:     if  $focus = +$  then
36:        $next2 \leftarrow 1$ 
37:     else
38:        $next2 \leftarrow 0$ 
39:     end if
40:   end while
41:    $i \leftarrow M, minlcp \leftarrow M, next1 \leftarrow 0$ 
42:   if  $focus = -$  then
43:      $next2 \leftarrow 1$ 
44:   else
45:      $next2 \leftarrow 0$ 
46:   end if
47:   while  $i > 1$  do
48:     while  $i > 1$  and  $GSA[i - next1] = GSA[i -$ 
     $next1 - 1]$  do
49:        $next1 \leftarrow next1 + 1$ 
50:     end while
51:     while  $focus = -$  and  $i > 1$  and  $GSA[i -$ 
     $next1 - next2] = GSA[i - next1 - next2 - 1]$  do
52:        $next2 \leftarrow next2 + 1$ 
53:     end while
54:     for  $j$  from  $next1 + next2$  to 0 do
55:        $minlcp \leftarrow \min\{minlcp, LCP[i - j]\}$ 
56:       if  $j \geq next2$  then
57:          $LCP_d^*[i - j] \leftarrow \max\{LCP_d^*[i -$ 
     $j], minlcp\}$ 
58:       end if
59:     end for
60:      $i \leftarrow i - next1 - 1$ 
61:      $minlcp \leftarrow M, next1 \leftarrow 0$ 
62:     if  $focus = -$  then
63:        $next2 \leftarrow 1$ 
64:     else
65:        $next2 \leftarrow 0$ 
66:     end if
67:   end while
68:   return  $LCP_d^*$ 
69: end function

```

---

$begin_i + |s_i| - 1$ . At the end of this, add the sampled unique substrings in  $\mathcal{U}'_i$  to the hash table described in Section 2.1.

In the following we show that the above greedy strategy obtains the smallest number of unique substrings that cover the same set of unique  $L$ -mers as  $\mathcal{U}_i$ , provided that each unique substring in  $\mathcal{U}_i$  occurs only one time (i.e., any  $u_{i,x} \in \mathcal{U}_i$  is not identical to another  $u_{i,y} \in \mathcal{U}_i$  if  $x \neq y$ ). As a result, the total number of unique substrings in  $\mathcal{U}' = \cup_{i=1}^m \mathcal{U}'_i$  included in CAMMIQ index is also as small as possible.

**Theorem 6.** *If  $u_{i,x} \in \mathcal{U}_i \neq u_{i,y} \in \mathcal{U}_i$  for  $x \neq y$ , then GreedySampling returns the smallest  $\mathcal{U}'_i$  such that if an  $L$ -mer includes some  $u_{i,x} \in \mathcal{U}_i$ , then it also includes at least one  $u_{i,x'} \in \mathcal{U}'_i$ .*

*Proof.* Consider  $\mathcal{U}'_i = \{u'_{i,1}(l'_1, r'_1), \dots, u'_{i,|\mathcal{U}'_i|}(l'_{|\mathcal{U}'_i|}, r'_{|\mathcal{U}'_i|})\}$  that GreedySampling returns; also consider an alternative sample  $\mathcal{U}''_i = \{u''_{i,1}(l''_1, r''_1), \dots, u''_{i,|\mathcal{U}''_i|}(l''_{|\mathcal{U}''_i|}, r''_{|\mathcal{U}''_i|})\}$  of  $\mathcal{U}_i$  that covers the same set of  $L$ -mers; assume both sets are sorted by the left indices ( $l'_1 < \dots < l'_{|\mathcal{U}'_i|}$ ;  $l''_1 < \dots < l''_{|\mathcal{U}''_i|}$ ). First, observe that  $l'_1 \geq l''_1$

---

**Algorithm 4** GreedySampling( $\mathcal{U}_i = \{u_{i,1}(l_1, r_1), u_{i,nu_i}(l_{nu_i}, r_{nu_i})\}$ ,  $L$ )

---

```

1:  $\mathcal{U}'_i \leftarrow \emptyset$ 
2:  $\text{cur} \leftarrow \text{begin}_i$  //  $\text{begin}_i$ : the beginning position of  $s_i$  in  $s$ 
3:  $u_{i,\text{last}}(l_{\text{last}}, r_{\text{last}}) \leftarrow \text{NIL}$ 
4: for  $x$  from 1 to  $nu_i$  do
5:   if  $u_{i,\text{last}} \neq \text{NIL}$  and  $r_x \geq \text{cur} + L$  then
6:      $\mathcal{U}'_i \leftarrow \mathcal{U}'_i \cup \{u_{i,\text{last}}(l_{\text{last}}, r_{\text{last}})\}$ 
7:      $\text{cur} \leftarrow l_{\text{last}} + 1$ 
8:   end if
9:    $u_{i,\text{last}}(l_{\text{last}}, r_{\text{last}}) \leftarrow u_{i,x}(l_x, r_x)$ 
10: end for
11: return  $\mathcal{U}'_i$ 

```

---

since  $u'_{i,1}$  is the rightmost unique substring in  $\mathcal{U}_i$  which is fully included in the leftmost unique  $L$ -mer. As a consequence we must have  $l'_2 \geq l'_1$  (and so on). Otherwise, if  $l'_1 = l'_2$  then  $u'_{i,2}(l'_2, r'_2)$  is not the rightmost unique substring in the next unique  $L$ -mer whose left index is greater than  $l'_1$ , and if  $l'_1 > l'_2$  then there is some unique  $L$ -mer not covered by any unique substring in  $\mathcal{U}'_i$ . Therefore  $|\mathcal{U}'_i| \leq |\mathcal{U}''_i|$  since there is an injection between the elements in  $\mathcal{U}'_i$  and  $\mathcal{U}''_i$  until reaching the last unique  $L$ -mer on  $s_i$ .  $\square$

**Corollary 7.** *If  $u_{i,x} \in \mathcal{U}_i \neq u_{i,y} \in \mathcal{U}_i$  for  $x \neq y$  on each  $s_i$ , then  $\mathcal{U}' = \cup_{i=1}^m \mathcal{U}'_i$  is the smallest set of unique substrings such that if an  $L$ -mer on  $s_i$  includes some  $u_{i,x} \in \mathcal{U}_i$ , then it also includes at least one  $u_{i,x'} \in \mathcal{U}'_i$ .*

We note that the greedy sampling strategy works in practice even if there are actually unique substrings occurring more than once in a given genome, meaning  $u_{i,x} = u_{i,y} \in \mathcal{U}_i$  for some  $y > x$ , leading to  $\mathcal{U}'_i$  (and thus  $\mathcal{U}'$ ) being close to optimality, since these unique substrings would constitute a very small proportion ( $\leq 0.1\%$ ) with the default minimum length  $L_{\min} = 26$  of unique substrings. This gives significantly smaller indices than alternative  $k$ -mer based tools, and results in an integer program with a small number of constraints.

We applied the above strategy to sample doubly-unique substrings in  $\mathcal{D}_i$ , to obtain the minimum size  $\mathcal{D}'_i$  for each genome  $s_i$  so that the aggregate set of doubly-unique substrings  $\mathcal{D}'$  is at most twice the optimal, provided that each doubly-unique substring occurs once in each of the corresponding genomes.

**Corollary 8.** *If  $d_{i,x} \in \mathcal{D}_i \neq d_{i,y} \in \mathcal{D}_i$  for  $x \neq y$  on each  $s_i$ , then  $\mathcal{D}' = \cup_{i=1}^m \mathcal{D}'_i$  is at most twice as large as the smallest set of doubly-unique substrings such that if an  $L$ -mer on  $s_i$  includes some  $d_{i,x} \in \mathcal{D}_i$ , then it also includes at least one  $d_{i,x'} \in \mathcal{D}'_i$ .*

## 5.4 Proof of Theorem 1

We begin with the following theorem from Weissman et al. [62] that bounds the L1 distance between the empirical distribution of a sequence of independent, identically distributed random variables and the true distribution.

**Theorem 9.** *Let  $P$  be a probability distribution on the set  $\mathcal{A} = \{1, \dots, a\}$ . Let  $X_1, X_2, \dots, X_n$  be i.i.d. random variables distributed according to  $P$ . Then, for any given  $\epsilon > 0$ ,*

$$\Pr[||P - \bar{P}||_1 \geq \epsilon] \leq (2^a - 2) \exp(-n\epsilon^2/2)$$

where  $\bar{P}$  is the empirical estimation of  $P$  defined as  $\bar{P}(i) = \frac{\sum_{j=1}^n \Delta(X_i=j)}{n}$ , where  $\Delta(e) = 1$  if and only if  $e$  is true and  $\Delta(e) = 0$  otherwise.

Now consider a collection of genomes  $\mathcal{A} = \{s_1, \dots, s_a\}$  with relative abundances  $p_1, \dots, p_a$  and the set  $\mathcal{Q} = \{r_1, \dots, r_n\}$  of  $n$  reads (i.e.,  $L$ -mers) sampled independently and uniformly at random from  $\mathcal{A}$  according to  $p_1, \dots, p_a$ . On each genome  $s_i$  let  $n_i^L$  denote the total number of  $L$ -mers and  $q_i = \frac{nu_i^L}{n_i^L}$  be the proportion of unique  $L$ -mers; then the probability of a read  $r_j \in \mathcal{Q}$  corresponds to a unique  $L$ -mer on  $s_i$  is  $\frac{p_i n_i^L}{\sum_{i'=1}^a p_{i'} n_{i'}^L} \cdot q_i$ , and the probability of a read  $r_j \in \mathcal{Q}$  does not correspond to any unique  $L$ -mers on  $s_i$  is  $\frac{p_i n_i^L}{\sum_{i'=1}^a p_{i'} n_{i'}^L} \cdot (1 - q_i)$ .

Therefore  $r_1, \dots, r_j$  are i.i.d. distributed according to  $(\frac{p_1 n_1^L}{\sum_{i'=1}^a p_{i'} n_{i'}^L} \cdot q_1, \dots, \frac{p_a n_a^L}{\sum_{i'=1}^a p_{i'} n_{i'}^L} \cdot q_a, \sum_{i=1}^a \frac{p_i n_i^L}{\sum_{i'=1}^a p_{i'} n_{i'}^L} \cdot (1 - q_i)) = (p'_1 q_1, \dots, p'_a q_a, \sum_{i=1}^a p'_i (1 - q_i))$  where the last term corresponds to the reads that are not unique to any  $s_i \in \mathcal{A}$ .

**Proof of Theorem 1.** Let  $c_i$  be the number of reads assigned to  $s_i$ . Also let  $p' = (p'_1, \dots, p'_a)$ . We set  $\hat{p} = (\hat{p}_1, \dots, \hat{p}_a)$ , by defining  $\hat{p}_i = \frac{c_i/q_i}{n}$  to be the predicted abundance of  $s_i$  based on the number of reads assigned to it. Consider  $P = (p'_1 q_1, \dots, p'_a q_a, \sum_{i=1}^a p'_i (1 - q_i))$  and  $\hat{P} = (\frac{c_1}{n}, \dots, \frac{c_a}{n}, 1 - \sum_{i=1}^a \frac{c_i}{n}) = (\hat{p}_1 q_1, \dots, \hat{p}_a q_a, 1 - \sum_{i=1}^a \frac{c_i}{n})$ ; by definition, we have  $\|P - \hat{P}\|_1 \geq \sum_{i=1}^a |p'_i - \hat{p}_i| q_i \geq q_{\min} \cdot \sum_{i=1}^a |p'_i - \hat{p}_i| = q_{\min} \cdot \|p' - \hat{p}\|_1$ . Then the following three theorem statements hold.

- (i) Given that  $n \geq \frac{2(a+1) + \ln \frac{1}{\delta}}{(p_{\min} q_{\min})^2}$ , we have

$$\begin{aligned} \Pr \left[ \|p' - \hat{p}\|_1 \geq p_{\min} \right] &= \Pr \left[ q_{\min} \|p' - \hat{p}\|_1 \geq p_{\min} q_{\min} \right] \\ &\leq \Pr \left[ \|P - \hat{P}\|_1 \geq p_{\min} q_{\min} \right] \\ &\leq 2^{a+1} \exp(-n(p_{\min} q_{\min})^2/2) \\ &\leq 2^{a+1} \exp\left(-\frac{2(a+1) + \ln \frac{1}{\delta}}{(p_{\min} q_{\min})^2} (p_{\min} q_{\min})^2/2\right) \\ &= \frac{2^{a+1}}{e^{a+1}} \delta \\ &\leq \delta. \end{aligned}$$

This implies that with probability  $\geq 1 - \delta$  the L1 distance between  $p'$  and  $\hat{p}$  is upper bounded by  $p_{\min}$ . As a result we have  $\hat{p}_i > 0$  for each  $\hat{p}_i$ , i.e.  $c_i \geq 0$ .

- (ii) The proof follows by simply replacing  $p_{\min}$  with  $\epsilon$  in the proof of (i).
- (iii) The proof follows by simply replacing  $p_{\min}$  with  $\sqrt{\frac{2(\ln \frac{1}{\delta} + (a+1))}{n q_{\min}^2}}$  in the proof of (i). Specifically,

$$\begin{aligned} \Pr \left[ \|p' - \hat{p}\|_1 \geq \sqrt{\frac{2(\ln \frac{1}{\delta} + (a+1))}{n q_{\min}^2}} \right] &= \Pr \left[ q_{\min} \|p' - \hat{p}\|_1 \geq q_{\min} \cdot \sqrt{\frac{2(\ln \frac{1}{\delta} + (a+1))}{n q_{\min}^2}} \right] \\ &\leq \Pr \left[ \|P - \hat{P}\|_1 \geq \sqrt{\frac{2(\ln \frac{1}{\delta} + (a+1))}{n}} \right] \\ &\leq 2^{a+1} \exp\left(-\frac{n}{2} \cdot \frac{2(\ln \frac{1}{\delta} + (a+1))}{n}\right) \\ &= \frac{2^{a+1}}{e^{a+1}} \delta \\ &\leq \delta. \end{aligned}$$

## 5.5 blastn analysis of *Salmonella* strains

The RNASeq reads for each cell studied in [58] are stored as separate read sets in the Sequence Read Archive (SRA) [63]. To reduce the number of reads that need to be downloaded or aligned, we used ReadFinder (<https://github.com/morgulis/ReadFinder>) to find any reads that could plausibly map to either the *Salmonella* strains *LT2* or *D23580* with permissive parameters allowing alignments that stray by as much as four diagonals from the main diagonal. ReadFinder uses similar methods as SRPRISM [64] to search SRA without the need to download the data and is more streamlined than SRPRISM for our purpose of finding candidate matching reads. We then used blastn [65] with word size 16 to find local (and ideally, global) alignments between the reads identified by ReadFinder and either *Salmonella* strain. This test is

much simpler than the CAMMiQ test because we did not consider alignments to any *Salmonella* strains other than the two actually used in the original wet lab experiment.

Most reads that align to either strain actually align to the two strains equally well. To decide which reads align strictly better to one strain or the other, we implemented an in-house script with the following rules. A blastn alignment is "high-quality" if it has length at least 70 (taking into account that the reads are of length 75), has identity percentage  $\geq 98$ . A read  $R$  maps better to the *LT2* strain if either:

- $R$  has a high-quality alignment to the *LT2* strain but  $R$  has no high-quality alignment to the *D23580* strain or
- $R$  has high-quality alignments to both strains and the best *LT2* alignment is at least as good as the best *D23580* alignment on i) length ii) identity percentage, iii) gaps and is better than the best D-strain alignment on at least one of the three Roman numeral criteria.

The criteria for mapping better to the D-strain are symmetric. The script reports counts of how many reads map align strictly better to each of the two strains. Although the data in [58] consist of paired reads, each mate pair was treated individually in the blastn analysis.

## 6 Supplementary Figures

Figure 5 presents 50 genomes in our species-level dataset with the least number of unique  $L$ -mers. In this graph each node represents one such genome; each edge connects two nodes if they share a doubly-unique substring. Solid black edges indicate a pair of nodes that share at least 30 doubly-unique substrings; the remaining edges in grey indicate node pairs with fewer number of shared doubly-unique substrings. Notice that there is a special node in the center, representing the union of all genomes not included in these 50-genome subset. Any node connected to this special node by a single edge, or by a path, is identifiable and quantifiable by CAMMiQ, provided that all edges in this path are black (22 of these 50 nodes are as such) or all nodes in this path have sufficient abundance. Note that 20 of the genomes here are connected to this special node by a black edge: these are the genomes that form the least-quantifiable-20 dataset in our experiments.



