

# A computational toolset for rapid identification of SARS-CoV-2, other viruses and microorganisms from sequencing data

Shifu Chen, Changshou He, Yingqiang Li, Zhicheng Li and Charles E Melançon III

Corresponding author: Shifu Chen, Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, 518055 Guangdong, China.  
Tel: +86 135 1042 7830; Fax: +86 755 8657 4950; E-mail: chen@haplox.com

## Abstract

In this paper, we present a toolset and related resources for rapid identification of viruses and microorganisms from short-read or long-read sequencing data. We present *fastv* as an ultra-fast tool to detect microbial sequences present in sequencing data, identify target microorganisms and visualize coverage of microbial genomes. This tool is based on the k-mer mapping and extension method. K-mer sets are generated by UniqueKMER, another tool provided in this toolset. UniqueKMER can generate complete sets of unique k-mers for each genome within a large set of viral or microbial genomes. For convenience, unique k-mers for microorganisms and common viruses that afflict humans have been generated and are provided with the tools. As a lightweight tool, *fastv* accepts FASTQ data as input and directly outputs the results in both HTML and JSON formats. Prior to the k-mer analysis, *fastv* automatically performs adapter trimming, quality pruning, base correction and other preprocessing to ensure the accuracy of k-mer analysis. Specifically, *fastv* provides built-in support for rapid severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) identification and typing. Experimental results showed that *fastv* achieved 100% sensitivity and 100% specificity for detecting SARS-CoV-2 from sequencing data; and can distinguish SARS-CoV-2 from SARS, Middle East respiratory syndrome and other coronaviruses. This toolset is available at: <https://github.com/OpenGene/fastv>.

**Key words:** SARS-CoV-2; viruses; microorganisms; identification; k-mer

## Introduction

The coronavirus disease-2019 (COVID-19) pandemic has spread to over 200 countries and territories, and it has made a terrible impact on lives and economies worldwide [1–3]. Based on

the current pandemic situation and many research reports [4], COVID-19 may continue to spread for a long period of time and may eventually become a flu-like seasonal outbreak [5]. Under these circumstances, it is important to develop new technologies for rapid detection of COVID-19.

**Shifu Chen** is a researcher at Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences. He also serves as chief technology officer of HaploX Biotechnology. He is the initiator of OpenGene projects and a contributor to many open source tools.

**Changshou He** is an engineer at department of bioinformatics, HaploX Biotechnology.

**Yingqiang Li** is an engineer at department of bioinformatics, HaploX Biotechnology.

**Zhicheng Li** is a researcher at Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences. His research interests lie mainly in imaging genomics.

**Charles E Melançon III** is a senior scientist at department of research and development, HaploX Biotechnology. His research interests lie mainly in next-generation sequencing and bioinformatics.

Submitted: 22 May 2020; Received (in revised form): 3 August 2020

© The Author(s) 2020. Published by Oxford University Press. All rights reserved. For Permissions, please email: journals.permissions@oup.com

Nucleic acid sequencing is a key technology for identifying and studying severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), the causative agent of COVID-19. Metagenomic and metatranscriptomic next-generation sequencing (mNGS) are powerful tools to study the genetic composition and function of microbial populations, as well as to analyze the relationship between microorganisms and their host or environment [6, 7]. Since the first clinical application of metagenomic sequencing (mNGS) for the diagnosis of leptospirosis in 2014 [8], mNGS has been used widely for the identification and diagnosis of new and rare pathogens.

Next generation sequencing (NGS) technology has also played an important role in COVID-19 diagnosis and research. The earliest COVID-19 case, which was initially diagnosed as pneumonia caused by an unknown pathogen, was identified due to the presence of SARS-like sequences found using mNGS [9]. The first complete genome of SARS-CoV-2 (GenBank: MN908947) reported on January 11, 2020 was assembled using NGS data [10]. Soon after, the whole-genome sequence of SARS-CoV-2 was also obtained by mNGS using the Oxford Nanopore platform supplemented with Sanger sequencing [11]. The rapid acquisition and publication of the SARS-CoV-2 genome was essential to design fluorescent polymerase chain reaction (PCR) probes for COVID-19 nucleic acid detection kits.

With the viral genome in hand, we can now explore the possibility of using mNGS directly as a detection method to determine whether SARS-CoV-2 RNA is present or absent in a sample. In theory, a simple and straightforward approach would be to first map sequencing reads obtained from the sample to the viral genome using common aligners such as Burrows-Wheeler aligner (BWA) [12] or Bowtie2 [13], and then to analyze the alignment results to determine the coverage of the viral genome and the number of properly mapped reads. However, in practice, such an alignment-based method is prone to problems stemming from both false positives and false negatives. On one hand, some viruses have genomes very similar to SARS-CoV-2, which can lead to false positive results. For example, the genome of bat coronavirus RaTG13 is sufficiently similar to SARS-CoV-2 (96% identity) to cause non-specific alignment [14]. On the other hand, in some cases the virus-specific reads obtained may not be abundant enough for unambiguous detection, which can lead to false negative results. Examples of such cases may be when the viral RNA is highly degraded, or when the sequencing library has been incompletely target enriched by multiple PCR [15] or hybrid capture [16]. In these scenarios, alignment-based methods may be not specific or sensitive enough. Such alignment-based methods are also computationally intensive and therefore not particularly fast or efficient.

In order to detect SARS-CoV-2 more quickly and accurately, we have developed an alignment-free method based on k-mer mapping and extension. The use of k-mer-based methods to analyze microbial sequencing data is not a new method. For example, researchers have developed k-mer based algorithms to analyze plasmid-derived sequence fragments in metagenomics data [17, 18]. Several k-mer based tools for species classification and sequencing read annotation have also been reported. For example, SPINGO [19] provides rapid species classification for microbial amplicon sequences based on k-mer mapping technology. Kraken2 [20], a very popular taxonomic classification system, is also based on k-mer matches. KrakenUniq [21], which is built on Kraken [22], can partially address problems due to false positives by using unique k-mer analysis. However, we are currently unaware of a fast, reliable and user-friendly k-mer-based tool for identifying nucleic acids from SARS-CoV-2 and

other viruses or microorganisms using sequencing data. This unmet need has led us to develop a new toolset and to provide corresponding k-mer resources.

Here, we present *fastv* and *UniqueKMER*, along with the precomputed unique k-mer resources. The *fastv* tool has three major functions: (i) to analyze which viral and/or microbial sequences are present in the sequencing data, (ii) to determine whether sequences from a specific virus or microorganism (e.g. SARS-CoV-2) can be found in the sequencing data and (iii) to analyze coverage of a specific viral or microbial genome by a set of sequencing data. All three of these methods rely on a unique k-mer set for each microorganism, making it essential to produce high-quality unique k-mer sets. We developed another tool, *UniqueKMER*, to generate a complete set of unique k-mers for each of a large set of microbial genomes. The unique k-mers can be filtered to remove the k-mers that can be mapped to a reference genome (e.g. the human genome). Generating unique k-mers for tens of thousands of viral and microbial species would typically require tremendous memory and computing resources. We have designed efficient algorithms to make this computation feasible on ordinary computing servers. *UniqueKMER* requires only 1 h to generate reference-filtered unique k-mers for about 12 000 viral genomes on an ordinary PC.

As a lightweight tool, *fastv* accepts FASTQ data as input and directly outputs the results in HTML and JSON formats. The HTML result is highly informative and provides interactive reports for manually reading, while the JSON result is structured such that it can easily be used by downstream analysis tools. Prior to the k-mer analysis, *fastv* automatically performs adapter trimming, quality pruning, base correction and other preprocessing to ensure the accuracy of k-mer analysis. These preprocessing features are derived from *fastp* [23], a popular quality control and filtering tool for NGS data previously developed by our group. The *fastv* tool is ultra-fast - it can process 10 M+ bases per second - and can complete the processing of a typical mNGS dataset in a few minutes.

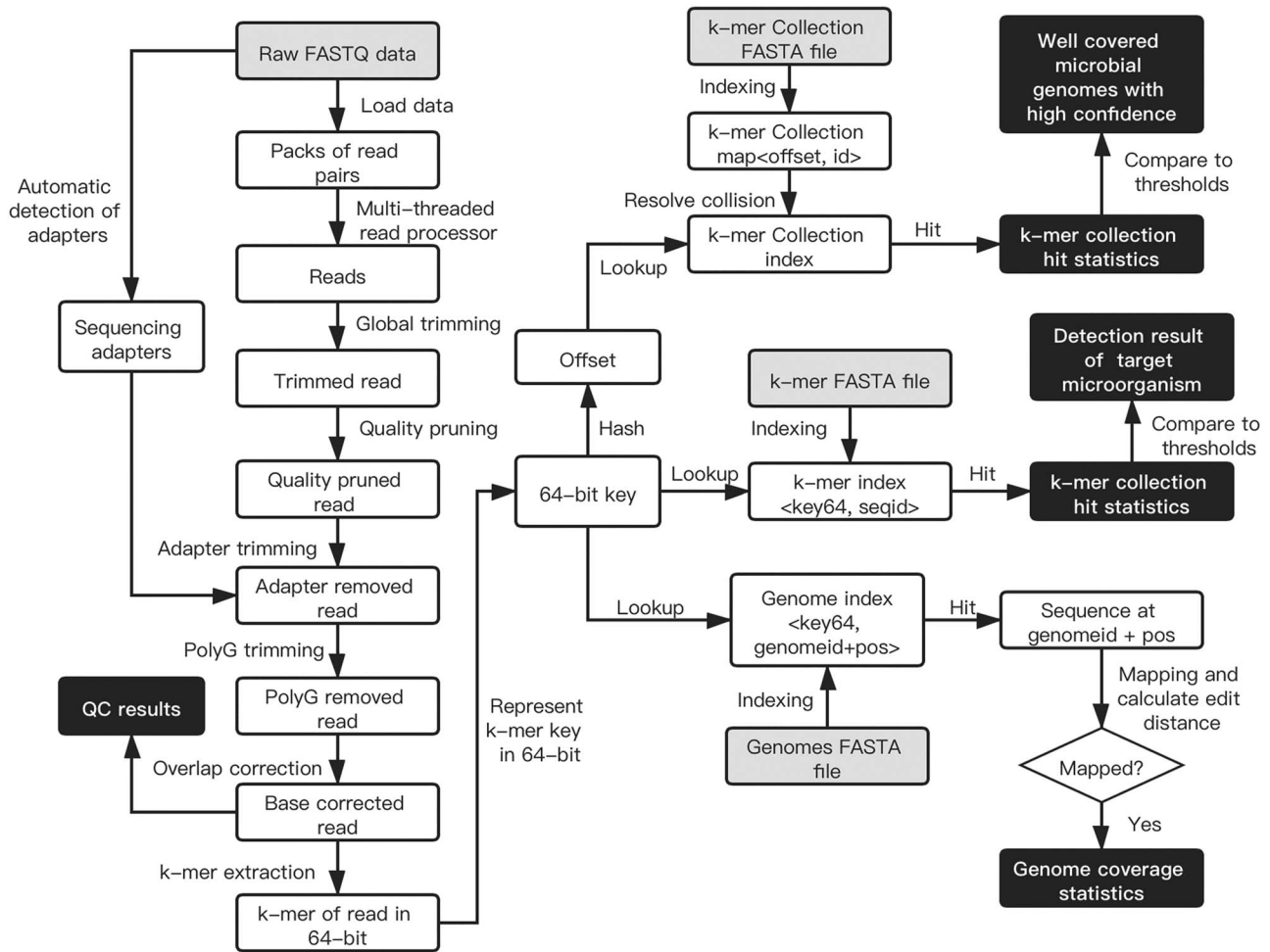
We conducted identification experiments on 27 samples positive for SARS-CoV-2 and 25 samples negative for SARS-CoV-2. The results showed that *fastv* achieved 100% sensitivity and 100% specificity; and that it can distinguish SARS-CoV-2 from SARS [24], middle East respiratory syndrome (MERS) [25] and other coronaviruses [26]. Although our original intention in developing these tools was to quickly identify SARS-CoV-2 from sequencing data, our tool can detect any target virus or microorganism for which a unique k-mer file is provided. We also conducted experiments using several other viral genomes such as Epstein-Barr virus (EBV), human papillomavirus (HPV) and hepatitis B virus (HBV). The results demonstrated that our tools perform well on a variety of viral genomic datasets.

## Material and methods

This section consists of three subsections: rapid identification of microorganisms from sequencing data, algorithms for efficiently generating unique k-mer sets and pregeneration of unique k-mer sets for common viruses and microorganisms.

### Fastv: rapid identification of microorganisms from sequencing data

*Fastv* is a highly optimized FASTQ scanner and k-mer mapper. Sequencing reads are first preprocessed to remove adapters and unqualified bases. Continuous k-mers are then computed for each read to be mapped to unique k-mer indexes. *Fastv* accepts



**Figure 1.** Overview of the fastv workflow. The items with grey backgrounds are input files, while the items with black backgrounds are results that will be output to HTML/JSON reports. An individual thread loads the FASTQ data to read packs (pack size = 1000). Multi-threaded read processors process data pack by pack. For each read or read pair, preprocessing is performed as in fastp. k-mers are extracted from the read and its reverse complement and are then converted to 64-bit keys. Each 64-bit key is used to search one of the three k-mer indexes built from the k-mer collection file, k-mer file, or genome(s) file.

input from any or all of the following three FASTA file types to generate corresponding k-mer indexes:

- Unique k-mer collection for a large set of viruses or microorganisms. This file contains a list of viruses or microorganisms along with their unique k-mers. The identifier of each FASTA entry represents the name of a viral or microbial genome, while its corresponding multiline sequences represent its unique k-mer keys. This file typically consists of all viruses or microorganisms that might be detected, with reference genomes available.
- Unique k-mer set for a specific virus or microorganism. This file contains a list of k-mer keys unique to a specific virus or microorganism. The sequence of each FASTA entry represents a k-mer key, while its corresponding FASTA identifier represents its position in the viral or microbial genome. This file typically represents the unique k-mer set for the virus or microorganism of interest, such as SARS-CoV-2.
- Genome sequences for a specific virus or microorganism. This file contains one or more reference genomes for the target virus or microorganism. Typically, multiple genomes represent different subtypes of the target virus or microorganism. For instance, if the target virus is HPV, the genomes

may comprise HPV-16, HPV-18, HPV-31, etc [27]. Due to the differences among different genome sequences, the coverage and mismatch rate of each genome will be different. This information can be used for microbial subtype identification.

Extraction of k-mers from reads and processing of k-mer indexes are independent processes. Figure 1 summarizes the fastv workflow.

#### Read preprocessing

To avoid generating erroneous k-mer keys, sequencing adapters and low-quality bases must be removed. Fastv utilizes the adapter cutting and quality pruning features from fastp, which we developed previously. Adapter sequences can either be auto-detected or specified from the command line. For paired-end sequencing data, overlapping regions are detected and incorrect bases in the overlapping regions are corrected. The algorithms for, and implementations of, these features can be found in the fastp publication [23]. For data generated by long-read platforms (e.g. PacBio or Oxford Nanopore Technologies (ONT)), long reads are segmented, generating multiple short reads.

### K-mer generation and representation

To accelerate k-mer generation and lookup, we use a unique 64-bit integer to represent each k-mer. A base (A/T/C/G) is represented by two bits, so a 64-bit integer can represent up to a 32-mer, which is sufficiently long for identifying a virus or microorganism. Any k-mer key that contains a degenerate base (i.e. N base) is ignored. This k-mer representation has been used widely in our previous works, such as GeneFuse [28]. A progressive method for k-mer calculation is applied to accelerate k-mer generation. The formula to compute the  $n$ th k-mer key can be denoted as:  $kmer(n) = ((kmer(n-1) \ll 2) + base2bit(n)) \& bit\_mask$ , where  $base2bit(n)$  is the 2-bit representation of a base (A/T/C/G), and  $bit\_mask$  is a 64-bit value determined by the length  $k$  of k-mer.

### K-mer collection scanning

The k-mer collection file can contain unique k-mers for tens of thousands of viral or microbial genomes, with each containing hundreds to thousands of k-mer keys. Therefore, there may be tens of millions of k-mer keys in total. It would be trivial to use a map<key, id> to index this data, but such an implementation would result in very slow access. Our approach was to build a hash function to hash the 64-bit key to a larger number (i.e.  $2^{30}$ ), which stores the index to the k-mer element. Because of the nature of the hash function, two different keys may have the same hash value, resulting in a hash collision [29]. However, as long as the hash function is sufficiently random, the probability of collision will follow a probability distribution that will result in only a small fraction of keys to have hash collisions. This kind of space-for-time approach results in a significant increase in efficiency and makes it feasible to detect tens of thousands of microorganisms at once. The algorithm to build such a k-mer collection index is briefly illustrated as Algorithm 1. It should be noted that, to obtain a higher running speed, we did not use mutex locks to synchronize the k-mer hit counting operation between multiple threads. Our results demonstrate that, while this might lead to slight instability in the results (less than 1 in  $10^4$ ), it does not have an impact on the overall quality of the results.

**Algorithm 1:** Generation of k-mer collection index

---

```

HASH_LENGTH = (1 < < 30)
index = array(HASH_LENGTH)
tmp_valid_kmer = array()
final_valid_kmer = array()
initialize(index, undefined)
for id in genomes_of_kmer_collection
  for seq in kmer_of_genome(id)
    key64 = seq_to_key64(seq)
    offset = hash(key64)
    if index[offset] != undefined and index[offset] != id
      index[offset] = collided
    else
      index[offset] = id
      tmp_valid_kmer.add(id, key64)
counter = 0
for id, key64 in tmp_valid_kmer
  offset = hash(key64)
  if index[offset] == id
    final_valid_kmer.add(id, key64)
    index[offset] = counter
    counter++

```

---

The hash function used in Algorithm 1 is a simple formula that can be calculated efficiently. We have used this hash function in previous work [30]. It utilizes the multiplication and bit manipulation of the key with several large prime numbers:

$$\begin{aligned} \text{hash}(\text{key64}) = & (1713137323 * \text{key64} + (\text{key64} \gg 12) \\ & * 7341234131 + (\text{key64} \gg 24) * 371371377) \& \\ & (\text{HASH\_LENGTH} - 1) \end{aligned}$$

### K-mer scanning for a specific virus or microorganism

Since the k-mer list of a specific virus or microorganism is typically small, it is trivial to implement k-mer scanning on such a short list. A simple unordered map is used to represent the index for such a k-mer set, with the values used as k-mer hit statistics. In contrast to k-mer collection scanning, where multiple threads share a single hit counting statistics array, each thread in the k-mer scanning operation uses its own statistics array. These arrays are subsequently merged to generate the overall statistics, making the k-mer scanning result stable and reproducible.

### Genome coverage statistics and subtyping

The genomes are indexed as a map, with its key as the same 64-bit k-mer key, and its value as a list of genome positions (GP). A genome position includes a genome ID and a position in that genome. For a given read, if one of its 64-bit k-mer keys is a hit to the genome index, the read will be mapped to the corresponding genome location. The edit distance [31] between the read and the genome sequence will be obtained; and if the edit distance is less than or equal to the threshold, a match will be recorded and the coverage of that genome will be updated. It is worth mentioning that some microbial genomes (i.e. EBV) have a large number of repeated sequences [32], which will result in some keys being hits to many different genome positions. In such cases, fastv divides the coverage and mismatch numbers of a read into multiple parts and distributes them to each position, producing smooth and uniform coverage. When multiple genomes are input, the coverage result will be sorted by the coverage. This feature can help identify a subtype of a specific virus or microorganism.

### Visualization

The k-mer scanning results of different inputs are visualized in a figure on a single HTML page. For k-mer scanning of a specific virus or microorganism, the result is simply plotted with the widely used library Plotly.js. For genome k-mer scanning results, we developed a much more efficient toolkit based on native browser utilities to illuminate genome coverage and mismatch ratios. Our highly optimized method can easily support the visualization of hundreds of genomes, which is impossible for most common plotting libraries. A demonstration of a fastv HTML report is shown in Figure 2.

### UniqueKMER: efficient unique k-mer generation for large datasets

Since the key features of fastv rely on unique k-mer mapping and extension, it is important to obtain high-quality unique k-mer sets for microorganisms of interest. Although a number of k-mer generation tools are currently available [33–35], none



are suitable for our application because we must both generate unique k-mers for tens of thousands of viruses and/or microorganisms and filter the k-mer keys based on the reference genome. For instance, KMC3 [35] is an excellent tool for handling k-mer generation for each genome, but is inconvenient and less efficient for generating unique k-mers. It also is limited in the number of genomes (255 for the current version) for which it can perform unique k-mer extraction, thus making it unsuitable for generating unique k-mers for thousands of genomes. These unmet needs have led us to develop UniqueKMER, a new unique k-mer generation tool. The workflow of UniqueKMER is briefly described in Figure 3.

As shown in Figure 3, the UniqueKMER workflow consists of two parts: unique k-mer generation and unique k-mer filtering. In the first part all k-mer keys are extracted, and keys that belong to more than one genome are removed as non-unique keys. In the second part, both the keys that exactly match the reference genome and the keys that can be partially mapped to reference genome are removed. Although this can be done with a common aligner such as BWA or Bowtie2, these are not ideal for partial mapping of short sequences. We developed an algorithm based on S-bit seeding and edit distance computation to address this problem, as briefly shown in Algorithm 2.

---

**Algorithm 2:** unique k-mer filtering by reference genome

---

```
// KMER length in bases, by default, kmer_len = 25
kmer_len = 25
unique_keys = compute_unique_keys(genomes, kmer_len)
// seed_len is much less than kmer_len
seed_len = min(14, kmer_len)
seed_bit_len = 2 * seed_len
bloom_filter_len = 1 < < seed_bit_len
bloom_filter = array(bloom_filter_len)
key_genome_pos = map < key, genome_pos_list >
// Initialize the bloom filter
for key64 in unique_keys
  for s_bit_key in tranverse_by_bits(key64):
    bloom_filter[s_bit_key] = true
// Index the reference genome
for s_bit_key, genome_pos in reference_genome
  if bloom_filter[s_bit_key] == true
    key_genome_pos[s_bit_key].add(genome_pos)
// Remove the unique KMER keys that can be mapped to
reference genome
for key64 in unique_keys
  for s_bit_key in tranverse_by_bits(key64)
    for genome_pos in key_genome_pos[s_bit_key]
      key_seq = to_sequence(key64)
      ref_seq = get_ref_sequence(reference_genome, genome_pos)
      ed = edit_distance(key_seq, ref_seq)
      if ed < THRESHOLD
        unique_keys.remove(key64)
```

---

**Pre-generation of unique k-mers for SARS-CoV-2 and other common viruses and microorganisms**

We pre-generated unique k-mers for two datasets. The first dataset is the National Center for Biotechnology Information (NCBI) viral genomes RefSeq database [36], which can be found at <https://ftp.ncbi.nlm.nih.gov/refseq/release/viral/>. The other pre-generated dataset is the NCBI human bacterial microbiome

RefSeq database [37], which can be found at [https://ftp.ncbi.nlm.nih.gov/genomes/HUMAN\\_MICROBIOM/Bacteria/](https://ftp.ncbi.nlm.nih.gov/genomes/HUMAN_MICROBIOM/Bacteria/). Because bacterial genomes often have multiple contigs, we concatenated contigs from a single bacterial genome by inserting 32 Ns between each contig, guaranteeing that no artifactual k-mer keys will be introduced unexpectedly. The generated resources can be found at the UniqueKMER repository (<https://github.com/OpenGene/UniqueKMER>).

SARS-CoV-2 is included in the viral genome list so its unique k-mer set was also generated. We selected 12 SARS-CoV-2 genomes from the GISAID database, which is a global initiative on sharing genomic data of influenza viruses and SARS-CoV-2. These genomes were not used for SARS-CoV-2 unique k-mer generation, but for coverage evaluation for genomes from different strains. Forster et al. recently conducted evolutionary analysis of 160 SARS-CoV-2 genomes [38]. They classified SARS-CoV-2 into three types (A, B and C) according to amino acid changes and determined the evolutionary relationships of these types. Based on this work, we selected two ancestral and derived genomes from each type, taking into consideration the date and location of the collected samples. Because SARS-CoV-2 has few mutations to date [39], the similarity between the genomes of different types is very high. Nevertheless, coverage sorting is able to identify the most closely related genome to the query sequence data. The selected SARS-CoV-2 genomes are available at the fastv repository (<https://github.com/OpenGene/fastv>) and will be updated periodically as new genomes are made available.

## Results

### SARS-CoV-2 identification

To evaluate the performance of fastv for SARS-CoV-2 identification, we conducted experiments on 27 samples that were positive for SARS-CoV-2 and 25 samples that were negative for SARS-CoV-2. The platforms used for sequencing these samples were diverse and included Illumina, Oxford Nanopore, BGI-Seq, Capillary (Sanger sequencing) and Ion Torrent. For comparison, we also conducted alignment-based identification of SARS-CoV-2 using the well-known aligners BWA and Bowtie2; and performed taxonomic classification with the widely used tools Kraken2 and KrakenUniq. For BWA, we applied the BWA-MEM algorithm, which performs local alignment. For Bowtie2, we applied the global alignment method. The results are shown in Table 1.

As shown in Table 1, fastv achieved 100% sensitivity and 100% specificity for all tested samples; and could distinguish SARS-CoV-2 from SARS, MERS and other coronaviruses. The ground truth results were collected from corresponding papers that reported the data. The fastv default cutoff setting was used for determining whether or not SARS-CoV-2 was present. The pipelines for alignment-based SARS-CoV-2 identification are described in Supplementary File 1 (BWA-MEM method) and Supplementary File 2 (Bowtie2 method). The two alignment-based pipelines both failed to identify a SARS-CoV-2 sample (SAMN14445407), which was previously enriched using multiplex PCR technology. They also incorrectly identified the bat coronavirus RaTG13 sample (SAMN14082201) as SARS-CoV-2 since genome of RaTG13 has about 96% similarity to the genome of SARS-CoV-2 [40]. Even if careful manual adjustment of the alignment parameters may result in a better result for this dataset, it is difficult to ensure that the adjusted parameters can achieve good results on other datasets. The fastv results are based on the default parameters determined before this

Table 1. Comparative performance of fastv, Kraken2, KrakenUniq and alignment-based methods for identification of SARS-CoV-2

BioSample accession number	Platform	Ground truth of the data	Fastv result	BWA-MEM result	BWA-MEM coverage (%)	Bowtie2 result	Bowtie2 coverage (%)	Kraken2 reads	Kraken2 rank	KrakenUniq reads	KrakenUniq rank
SAMN14422484	Ion Torrent	SARS-CoV-2	POSITIVE	POSITIVE	56.36	POSITIVE	30.49	215	6	263	7
SAMN14381187	ONT	SARS-CoV-2	POSITIVE	POSITIVE	29.35	POSITIVE	27.31	38 856	1	38 952	1
SAMN14380341	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	99.83	POSITIVE	99.83	485 029	1	515 516	1
SAMN14341640	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	99.98	POSITIVE	100	3 432 594	1	3 782 498	1
SAMN14341639	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	99.85	POSITIVE	99.81	110 153	1	121 445	1
SAMN14308026	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	100.00	POSITIVE	100	1 082 910	1	1 085 843	1
SAMN14332760	BGISEQ	SARS-CoV-2	POSITIVE	POSITIVE	36.20	POSITIVE	33.9	190	1	203	1
SAMN14308029	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	99.41	POSITIVE	99.35	4731	1	5947	1
SAMN14180202	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	100.00	POSITIVE	99.95	551 135	1	564 321	1
SAMN14154205	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	99.91	POSITIVE	99.78	1 005 299	1	1 006 180	1
SAMN14154204	ONT	SARS-CoV-2	POSITIVE	POSITIVE	100.00	POSITIVE	99.83	78 590	1	80 411	1
SAMN14154203	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	48.24	POSITIVE	18.9	37 364	1	37 399	1
SAMN14154202	ONT	SARS-CoV-2	POSITIVE	POSITIVE	100.00	POSITIVE	99.81	29 524	1	30 222	1
SAMN14154201	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	44.49	POSITIVE	14.99	43 163	1	43 170	1
SAMN14154200	ONT	SARS-CoV-2	POSITIVE	POSITIVE	100.00	POSITIVE	99.98	468 580	1	480 655	1
SAMN14154199	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	99.88	POSITIVE	99.76	29 685	1	33 204	1
SAMN14154198	ONT	SARS-CoV-2	POSITIVE	POSITIVE	100.00	POSITIVE	99.6	36 364	1	37 710	1
SAMN14082199	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	97.60	POSITIVE	97.43	1219	41	1222	9
SAMN14082197	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	99.96	POSITIVE	99.92	6076	12	6088	4
SAMN14082196	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	92.57	POSITIVE	92.13	8002	21	801	10
SAMN14082200	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	99.91	POSITIVE	99.8	2603	130	2609	24
SAMN13922059	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	99.44	POSITIVE	99.97	61 904	16	62 399	6
SAMN13872787	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	100.00	POSITIVE	99.87	14 019	1	14 071	2
SAMN13872786	Illumina	SARS-CoV-2	POSITIVE	POSITIVE	100.00	POSITIVE	99.99	55 170	1	55 452	2
SAMN13871323	ONT	SARS-CoV-2	POSITIVE	POSITIVE	99.78	POSITIVE	95.85	703	5	724	4
SAMN13898864	ONT	SARS-CoV-2	POSITIVE	POSITIVE	100.00	POSITIVE	99.96	231 846	1	243 212	1
SAMN14445407	BGISEQ	SARS-CoV-2	POSITIVE	NEGATIVE	13.97	NEGATIVE	13.84	63 602	1	1 306 752	1
SAMN14086238	Illumina	Bat Coronavirus	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMN14086235	Illumina	Bat Coronavirus	NEGATIVE	NEGATIVE	0.50	NEGATIVE	0.00	1	2375	1	138
SAMN14086234	Illumina	Bat Coronavirus	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMN14086233	Illumina	Bat Coronavirus	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMN14086230	Illumina	Bat Coronavirus	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMN14082201	Illumina	Bat coronavirus RaTG13	NEGATIVE	POSITIVE	96.04	POSITIVE	63.05	812	6	848	12
SAMN02688745	Illumina	MERS-CoV	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMN02688792	Illumina	MERS-CoV	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMN02688873	Illumina	MERS-CoV	NEGATIVE	NEGATIVE	0.12	NEGATIVE	0.00	NA	NA	NA	NA
SAMN02402894	Illumina	SARS-CoV	NEGATIVE	NEGATIVE	17.71	NEGATIVE	8.44	910	4	6	8
SAMN02402954	Illumina	SARS-CoV	NEGATIVE	NEGATIVE	17.50	NEGATIVE	6.74	365	2	2	5
SAMN02402960	Illumina	SARS-CoV	NEGATIVE	NEGATIVE	5.57	NEGATIVE	1.99	NA	NA	NA	NA

Continued

Table 1. Continued.

BioSample accession number	Platform	Ground truth of the data	Fastv result	BWA-MEM result	BWA-MEM coverage (%)	Bowtie2 result	Bowtie2 coverage (%)	Kraken2 reads	Kraken2 rank	KrakenUniq reads	KrakenUniq rank
SAMEA5841282	ONT	Human coronavirus 229E	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMEA5841281	ONT	Human coronavirus 229E	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMEA5841278	ONT	Human coronavirus 229E	NEGATIVE	NEGATIVE	0.13	NEGATIVE	0.00	NA	NA	NA	NA
SAMEA3931206	Illumina	Human coronavirus OC43	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMEA3931205	Illumina	Human coronavirus NL63	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMEA3931204	Illumina	Human coronavirus HKU1	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMEA3931214	Illumina	Human respirovirus 1	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMEA3931208	Illumina	Human bocavirus 1	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMEA3931225	Illumina	Respiratory syncytial virus	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMN03386977	Illumina	HKU5 Coronavirus	NEGATIVE	NEGATIVE	12.44	NEGATIVE	8.36	18 166	4	26 542	6
SAMN03386974	Illumina	Bat coronavirus HKU3	NEGATIVE	NEGATIVE	0.63	NEGATIVE	9.38	132 098	2	237 706	2
SAMN12113789	Capillary	Influenza A Virus HKU11	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA
SAMN12113790	Capillary	Influenza A Virus HKU12	NEGATIVE	NEGATIVE	0.00	NEGATIVE	0.00	NA	NA	NA	NA

Positive/Negative: SARS-CoV-2 was/was not identified in the sample of the row using the method of the column. Coverage (%) × 1 coverage percentage of SARS-CoV-2 in BWA and Bowtie2 alignment results. Kraken2/KrakenUniq reads: number of SARS-CoV-2 reads in Kraken2 or KrakenUniq results. Kraken2/KrakenUniq rank: read number rank of SARS-CoV-2 in Kraken2 or KrakenUniq results. The possible false positives and false negatives introduced by alignment-based methods (coverage cutoff=20%) are marked in red and purple, respectively. Possible false positives and false negatives in Kraken2/KrakenUniq results are not shown due to high variability in results using various cutoff thresholds.

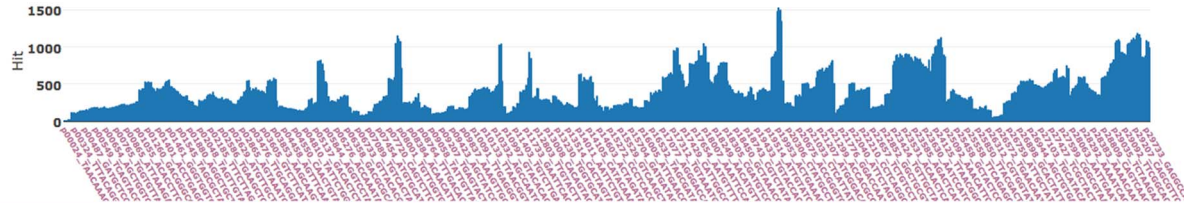
# fastv report

Created by [fastv](#) v0.7.0, an ultra-fast tool for fast identification of SARS-CoV-2, other viruses, and other microorganisms from sequencing data

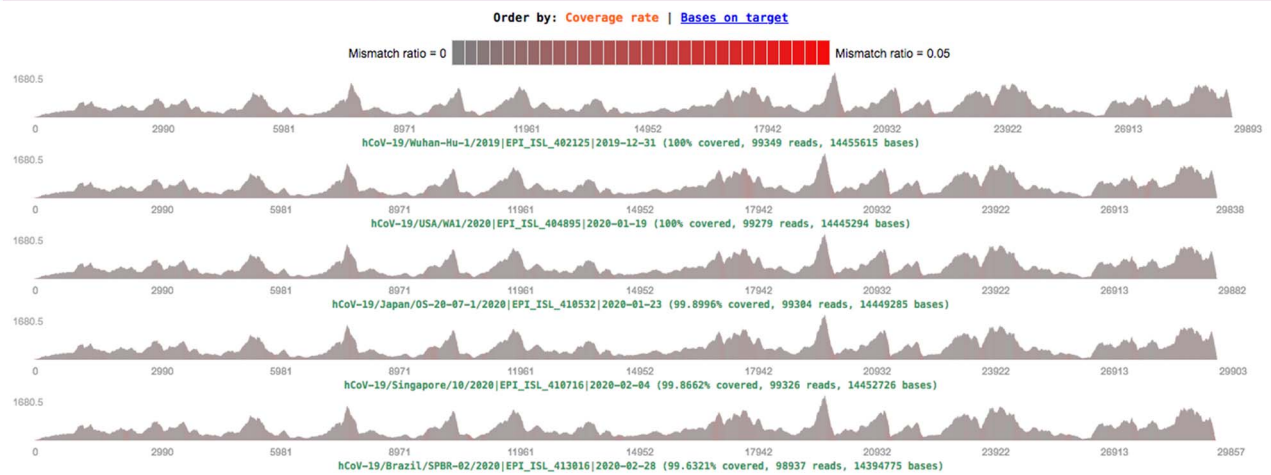
**Detection result for target unique k-mer file: [data/SARS-CoV-2.kmer.fa](#) (click to show/hide)**

Detection result for target k-mer file:	<b>POSITIVE</b>
Mean depth of k-mer coverage:	428.997238
Threshold to be positive:	0.100000

Unique k-mer hits (724 k-mer keys)



**Genome coverages for file: [data/SARS-CoV-2.genomes.fa](#) (click to show/hide)**



**Detection result for k-mer collection file: [/Users/shifu/data/virus/allkmers.fasta.gz](#) (click to show/hide)**

Genome	Coverage	Median depth	Mean depth	Remark
NC_045512.2 Wuhan seafood market pneumonia virus isolate Wuhan-Hu-1, complete genome (709 k-mer keys)	100%	364	429.298	SARS-CoV-2
NC_001422.1 Coliphage phi-X174, complete genome (488 k-mer keys)	99.5902%	15489	15294.5	Illumina PhiX control library

Figure 2. Fastv HTML report demonstration. The result for targeted k-mer hits is visualized using Plotly.js, whereas the result for genome coverage is visualized by a custom toolkit we developed. The k-mer collection scanning result shows that the data contains sequences of two microbial genomes. One is phi-X174, which is actually introduced by the Illumina PhiX control library, and the other is SARS-CoV-2. The genome coverage statistics show that SARS-CoV-2 most closely matches strain Wuhan-Hu-1. The red marks indicate the regions with a high mismatch ratio.

experiment. Therefore results obtained using fastv are more robust and reliable than those obtained using the BWA and Bowtie alignment-based methods.

As shown in the last four columns in Table 1, it is difficult to find a cutoff that reliably identifies SARS-CoV-2 using Kraken2 and KrakenUniq. Some positive samples had very few reads matching to SARS-CoV-2 (e.g. SAMN14332760 and SAMN14422484), whereas some negative samples had many reads matching to SARS-CoV-2 (e.g. SAMN03386977 and SAMN03386974). The SARS-CoV-2 read number ranks obtained with Kraken2 and KrakenUniq are also not good indicators of correct positive or negative SARS-CoV-2 identification. For example, a positive sample SAMN14082200 has a Kraken2 rank of 130 and a KrakenUniq rank of 24, whereas the Kraken2 and KrakenUniq ranks of a negative sample (SAMN03386974) are both 2.

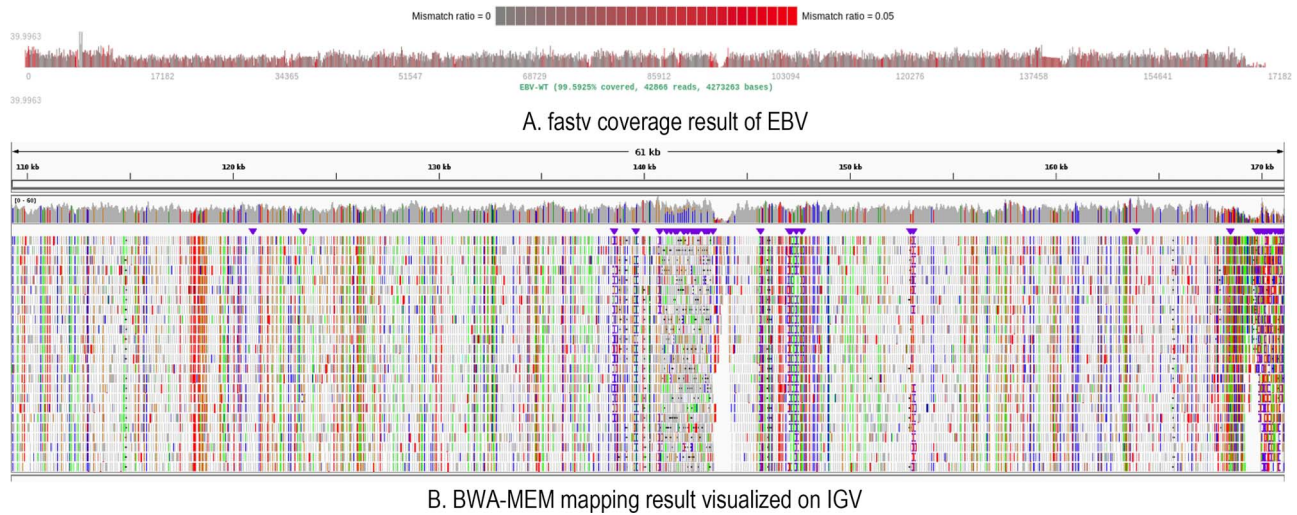
These results indicate that fastv achieved better performance for SARS-CoV-2 identification than alignment-based pipelines and taxonomic classifiers. Considering that it can also provide coverage and subtyping information, fastv is more suitable for detecting SARS-CoV-2 from sequencing data.

### Identification of other viruses

We also conducted experiments using data from other viruses and microorganisms. Epstein-Barr virus (EBV) has a long repetitive region within its genome, which often causes difficulties for other k-mer-based algorithms because these algorithms map a k-mer key to the location where it first appears. But our optimized k-mer mapping algorithm distributes hits corresponding to a k-mer key to all places where it appears, resulting in







**Figure 4.** EBV identification using fastv. The EBV genome has a large repetitive region between 12 kb and 35 kb. The coverage result generated by fastv (A) shows even coverage of this region and is similar to the IGV visualization of the BWA alignment result (B). The data used in this experiment was whole genome sequencing of an EBV-positive sample, downloaded from NCBI Sequence Read Archive (SRA accession: ERR1293949).



**Figure 5.** HBV subtyping using fastv. Eight HBV subtypes (HBV-A through HBV-H) are included in the genome list. The subtyping result is HBV-C, which is 99.69% covered. The HBV genome has many conserved regions where the genomes of different subtypes are very similar. This results in partial coverage of other HBV subtypes.

genes encoding hemagglutinin (HA gene) and neuraminidase proteins (NA gene) for over 38,000 influenza A strains from the Influenza Research Database [45] and constructed k-mer collection files for HA and NA genes. These were then used to identify the hemagglutinin and neuraminidase subtypes from a test dataset of 25 samples of different subtypes (H1N1, H2N2, H3N2, H5N1, H7N9 and H9N2). The results showed 24 samples were successfully classified according to their hemagglutinin and neuraminidase subtypes. The only sample (SRA accession: SRR5413408) for which hemagglutinin and neuraminidase subtypes could not be determined contained many amplicons with low coverage. The results show that fastv can be used for influenza A virus subtyping with high-quality sequencing data.

It should be noted that although fastv is effective at identifying major influenza A subtypes, it may not be suitable for distinguishing clades and sub-clades, since many sub-clades have an insufficient number of unique k-mers due to the high genomic similarity of different influenza A virus strains.

#### Identification of pathogen from mNGS data without set a target virus or microorganism

Fastv can also be used to quickly identify microorganisms and/or viruses in sequenced samples where the pathogen is unknown. To use this function, a k-mer collection file containing many possible viruses and microorganisms should first be prepared.

For user convenience we have pregenerated a k-mer collection file which from genomes of all human-associated viruses and bacteria with a reference genome provided in NCBI RefSeq database. This pregenerated k-mer collection file can be downloaded from the fastv repository. After scanning the FASTQ data, fastv will report the k-mer coverage for each microbial genome with valid hits. This information can be used to identify the pathogen. We evaluated seven mNGS datasets (SRA accessions: SRP006887, SRP006881, SRP000376, SRP007321, ERS4389819, SRP000657 and SRP004485), which were generated by Illumina and LS454 sequencers. All pathogens were correctly detected, with the k-mer coverage ranging from 15.19% to 99.94%.

## Discussion

In summary, we describe a new tool, fastv, for rapid identification of viruses and microorganisms from sequencing data. This tool is based on the k-mer mapping and extension method and relies on high-quality unique k-mers. We also describe a new tool, UniqueKMER, to generate such high-quality unique k-mer sets for a large collection of viruses and microorganisms. Experimental results show that with the k-mers generated by UniqueKMER, fastv is able to detect SARS-CoV-2 with 100% sensitivity and 100% specificity.

Because of the rapid and unpredictable spread of COVID-19, it is important to develop inexpensive, rapid and reliable methods for identification of its causative agent, SARS-CoV-2. Next-generation sequencing-based methods are suitable for SARS-CoV-2 detection and offer some advantages over other detection methods. Therefore, computational tools that can rapidly and reliably identify SARS-CoV-2 from sequencing data will be valuable to the research community. Fastv can also output on-target (e.g. SARS-CoV-2) clean reads to individual FASTQ files, which can be input to downstream analysis pipelines. For example, genome assembly with the on-target clean reads will be simpler and faster. The on-target reads can also be input to database search utilities like BLAST [46].

Although our original intention in developing fastv and UniqueKMER was to quickly identify SARS-CoV-2 from sequencing data, these tools can be used more generally to detect any target virus or microorganism for which unique k-mer files are provided. The results of our experiments with EBV, HPV and HBV, sequencing data demonstrate the general applicability of our tools. Because fastv can rapidly scan tens of thousands of genomes, it is a powerful tool for identifying pathogens from mNGS data. Fastv can also be used for rapid identification of viral and bacterial pathogens from a biological sample. We will continue to update our resource library so that researchers can directly use the pregenerated high-quality unique k-mer files for mNGS data analysis. Currently, fastv sorts the results according to the genome coverage and median k-mer hits without considering the pathogenicity of each virus or microorganism. In future work, we will incorporate data from microbial pathogen databases such as the Database for Reference Grade Microbial Sequences (FDA-ARGOS) [47], to provide further assist in pathogen identification. It should be noted that fastv is unable to identify novel pathogens that are previously unreported. As sequencing technology continues to progress, we will update this toolset as needed to ensure its suitability for pathogen identification using current sequencing technologies. K-mer resources will be updated at least once per quarter to include genomes of new viruses and microorganisms.

## Availability

As part of the OpenGene projects, fastv and UniqueKMER are open-sourced through the MIT license. Fastv is available at <https://github.com/OpenGene/fastv>, and UniqueKMER is available at <https://github.com/OpenGene/UniqueKMER>. The precomputed unique k-mer resources are also provided in these repositories.

### Key Points

- This tool presents a new tool fastv for rapid identification of SARS-Cov-2, other viruses and microorganisms.
- Another tool UniqueKMER is presented for generation of high-quality unique k-mers.
- Unique k-mer resources for tens of thousands of viruses and microorganisms have been precomputed and uploaded to the fastv repository.

## Supplementary data

Supplementary data are available online at <https://academic.oup.com/bib>.

## Acknowledgement

The authors would like to thank the user community for testing these tools and reporting bugs.

## Conflict of interest

The authors declare that there is no conflict of interest.

## Funding

The Shenzhen Science and Technology Program of China (JCYJ20170818160306270); project of Bureau of Industry and Information Technology of Shenzhen (Grant No. 20170922151538732); Shenzhen Science and Technology Innovation Committee Technical Research Project (Grant No. JSGG20180703164202084); and the projects of Development and Reform Commission of Shenzhen Municipality (Grant No. XMHT20190104006 and XMHT20200104013).

## References

1. Mahase E. Coronavirus covid-19 has killed more people than SARS and MERS combined, despite lower case fatality rate. *BMJ* 2020;**368**:m641.
2. McKibbin WJ, Fernando R. *The Global Macroeconomic Impacts of COVID-19: Seven Scenarios*, 2020, Australian National University, Australia. [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=3547729](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3547729).
3. Fernandes N. *Economic Effects of Coronavirus Outbreak (COVID-19) on the World Economy*, Available at SSRN 3557504, 2020. University of Navarra, Spain. [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=3557504](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3557504).
4. Gates B. Responding to Covid-19—a once-in-a-century pandemic? *N Eng J Med* 2020;**382**:1677–9.
5. Sajadi MM, Habibzadeh P, Vintzileos A, et al. *Temperature and Latitude Analysis to Predict Potential Spread and Seasonality for COVID-19*, Available at SSRN 3550308, 2020. University of Maryland, The United States. [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=3550308](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3550308).

6. Simmonds P, Adams MJ, Benkó M, et al. Consensus statement: virus taxonomy in the age of metagenomics. *Nat Rev Microbiol* 2017;15:161–8.
7. Mokili JL, Rohwer F, Dutilh BE. Metagenomics and future perspectives in virus discovery. *Curr Opin Virol* 2012;2:63–77.
8. Wilson MR, Naccache SN, Samayoa E, et al. Actionable diagnosis of neuroleptospirosis by next-generation sequencing. *N Engl J Med* 2014;370:2408–17.
9. Zhu N, Zhang D, Wang W, et al. A novel coronavirus from patients with pneumonia in China, 2019. *N Engl J Med* 2020;382:727–33.
10. Wu F, Zhao S, Yu B, et al. A new coronavirus associated with human respiratory disease in China. *Nature* 2020;579:265–9.
11. Chan JF-W, Yuan S, Kok K-H, et al. A familial cluster of pneumonia associated with the 2019 novel coronavirus indicating person-to-person transmission: a study of a family cluster. *Lancet* 2020;395:514–23.
12. Li H, Durbin R. Fast and accurate short read alignment with burrows-wheeler transform. *Bioinformatics* 2009;25:1754–60.
13. Langmead B, Salzberg SL. Fast gapped-read alignment with bowtie 2. *Nat Methods* 2012;9:357–9.
14. Zhang YZ, Holmes EC. A genomic perspective on the origin and emergence of SARS-CoV-2. *Cell* 2020;181:223–7.
15. Lundberg DS, Yourstone S, Mieczkowski P, et al. Practical innovations for high-throughput amplicon sequencing. *Nat Methods* 2013;10:999.
16. Duncavage EJ, Magrini V, Becker N, et al. Hybrid capture and next-generation sequencing identify viral integration sites from formalin-fixed, paraffin-embedded tissue. *J Mol Diagn* 2011;13:325–33.
17. Zhou F, Xu Y. cBar: a computer program to distinguish plasmid-derived from chromosome-derived sequence fragments in metagenomics data. *Bioinformatics* 2010;26:2051–2.
18. Zhou F, Olman V, Xu Y. barcodes for genomes and applications. *BMC Bioinformatics* 2008;9:1–11.
19. Allard G, Ryan FJ, Jeffery IB, et al. SPINGO: a rapid species-classifier for microbial amplicon sequences. *BMC Bioinformatics* 2015;16:324.
20. Wood DE, Lu J, Langmead B. Improved metagenomic analysis with kraken 2. *Genome Biol* 2019;20:257.
21. Breitwieser F, Baker D, Salzberg SL. KrakenUniq: confident and fast metagenomics classification using unique k-mer counts. *Genome Biol* 2018;19:1–10.
22. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biol* 2014;15:1–12.
23. Chen S, Zhou Y, Chen Y, et al. Fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics* 2018;34:i884–90.
24. Ksiazek TG, Erdman D, Goldsmith CS, et al. A novel coronavirus associated with severe acute respiratory syndrome. *N Engl J Med* 2003;348:1953–66.
25. Assiri A, McGeer A, Perl TM, et al. Hospital outbreak of Middle East respiratory syndrome coronavirus. *N Engl J Med* 2013;369:407–16.
26. Fan Y, Zhao K, Shi ZL, et al. Bat coronaviruses in China. *Viruses* 2019;11:210.
27. Speich N, Schmitt C, Bollmann R, et al. Human papillomavirus (HPV) study of 2916 cytological samples by PCR and DNA sequencing: genotype spectrum of patients from the west German area. *J Med Microbiol* 2004;53:125–8.
28. Chen S, Liu M, Huang T, et al. GeneFuse: detection and visualization of target gene fusions from DNA sequencing data. *Int J Biol Sci* 2018;14:843–8.
29. Liang RH, Mo T, Dong W, et al. Theoretical and experimental assessment of degenerate primer tagging in ultra-deep applications of next-generation sequencing. *Nucleic Acids Res* 2014;42:e98.
30. Chen S, Huang T, Wen T, et al. MutScan: fast detection and visualization of target mutations by scanning FASTQ data. *BMC Bioinformatics* 2018;19:16.
31. Gao X, Xiao B, Tao D, et al. A survey of graph edit distance. *Pattern Analysis Appl* 2009;13:113–29.
32. Falk K, Gratama J, Rowe M, et al. The role of repetitive DNA sequences in the size variation of Epstein–Barr virus (EBV) nuclear antigens, and the identification of different EBV isolates using RFLP and PCR analysis. *J Gen Virol* 1995;76:779–90.
33. Bose SM, Lalapura VS, Saravanan S, et al. k-core: Hardware Accelerator for k-mer Generation and Counting used in Computational Genomics. In: 2019 32nd International Conference on VLSI Design and 2019 18th International Conference on Embedded Systems (VLSID). 2019, p. 347–52. IEEE.
34. Pan T, Flick P, Jain C, et al. Kmerind: a flexible parallel library for k-mer indexing of biological sequences on distributed memory systems. *IEEE/ACM Trans Comput Biol Bioinform* 2017;16:1117–31.
35. Kokot M, Dlugosz M, Deorowicz S. KMC 3: counting and manipulating k-mer statistics. *Bioinformatics* 2017;33:2759–61.
36. Brister JR, Ako-Adjei D, Bao Y, et al. NCBI viral genomes resource. *Nucleic Acids Res* 2015;43:D571–7.
37. Pruitt KD, Tatusova T, Maglott DR. NCBI reference sequences (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res* 2007;35:D61–5.
38. Andersen KG, Rambaut A, Lipkin WI, et al. The proximal origin of SARS-CoV-2. *Nat Med* 2020;26:450–2.
39. Benvenuto D, Giovanetti M, Ciccozzi A, et al. The 2019-new coronavirus epidemic: evidence for virus evolution. *J Med Virol* 2020;92:455–9.
40. Zhang T, Wu Q, Zhang Z. Probable pangolin origin of SARS-CoV-2 associated with the COVID-19 outbreak. *Curr Biol* 2020;30:1346–51.
41. Suwannakarn K, Payungporn S, Chieochansin T, et al. Typing (a/B) and subtyping (H1/H3/H5) of influenza A viruses by multiplex real-time RT-PCR assays. *J Virol Methods* 2008;152:25–31.
42. Zou S, Han J, Wen L, et al. Human influenza A virus (H5N1) detection by a novel multiplex PCR typing method. *J Clin Microbiol* 2007;45:1889–92.
43. Ryabinin VA, Kostina EV, Maksakova GA, et al. Universal oligonucleotide microarray for sub-typing of influenza A virus. *PLoS One* 2011;6:e17529.
44. Zhao J, Ragupathy V, Liu J, et al. Nanomicroarray and multiplex next-generation sequencing for simultaneous identification and characterization of influenza viruses. *Emerg Infect Dis* 2015;21:400.
45. Zhang Y, Aevermann BD, Anderson TK, et al. Influenza research database: an integrated bioinformatics resource for influenza virus research. *Nucleic Acids Res* 2017;45:D466–74.
46. Johnson M, Zaretskaya I, Raytselis Y, et al. NCBI BLAST: a better web interface. *Nucleic Acids Res* 2008;36:W5–9.
47. Sichtig H, Minogue T, Yan Y, et al. FDA-ARGOS is a database with public quality-controlled reference genomes for diagnostic use and regulatory science. *Nat Commun* 2019;10:1–13.