

METHOD

Open Access



Matchtigs: minimum plain text representation of k -mer sets

Sebastian Schmidt^{1*}, Shahbaz Khan^{2*}, Jarno N. Alanko^{1,3}, Giulio E. Pibiri^{4,5} and Alexandru I. Tomescu^{1*}

*Correspondence:
sebastian.schmidt@helsinki.fi;
shahbaz.khan@cs.iitr.ac.in;
alexandru.tomescu@helsinki.fi

¹ Department of Computer Science, University of Helsinki, Helsinki, Finland

² Department of Computer Science and Engineering, Indian Institute of Technology Roorkee, Roorkee, India

³ Faculty of Computer Science, Dalhousie University, Halifax, Canada

⁴ Department of Environmental Sciences, Informatics and Statistics, Ca' Foscari University of Venice, Venice, Italy

⁵ ISTI-CNR, Pisa, Italy

Abstract

We propose a polynomial algorithm computing a *minimum* plain-text representation of k -mer sets, as well as an efficient near-minimum greedy heuristic. When compressing read sets of large model organisms or bacterial pangenomes, with only a minor runtime increase, we shrink the representation by up to 59% over unitigs and 26% over previous work. Additionally, the number of strings is decreased by up to 97% over unitigs and 90% over previous work. Finally, a small representation has advantages in downstream applications, as it speeds up SSHash-Lite queries by up to 4.26× over unitigs and 2.10× over previous work.

Keywords: k -mer sets, Plain text compression, Graph algorithm, Sequence analysis, Genomic sequences, Minimum-cost flow, Chinese postman problem

Background

Motivation

The field of k -mer-based methods has seen a surge of publications in the last years. Examples include alignment-free sequence comparison [1–3], variant calling and genotyping [4–8], transcript abundance estimation [9], metagenomic classification [10–13], abundance profile inference [14], indexing of variation graphs [15, 16], estimating the similarity between metagenomic datasets [17], species identification [18, 19] and sequence alignment to de Bruijn graphs [20–23]. All these methods are based mainly on k -mer sets, i.e. on the existence or non-existence of k -mers. They ignore further information like for example predecessor and successor relations between k -mers which are represented by the topology of a de Bruijn graph [24–26].

On the other hand, many classical methods such as genome assemblers [26–35] and related algorithms [36–38], are based on de Bruijn graphs and their topology. To increase the efficiency of these methods, the graphs are usually compacted by contracting all paths where all inner nodes have in- and outdegree one. These paths are commonly known as *unitigs*, and their first usage can be traced back to [39]. Since unitigs contain no branches in their inner nodes, they do not alter the topology of the graph,



© The Author(s) 2023. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

and in turn enable the exact same set of analyses. There are highly engineered solutions available to compute a compacted de Bruijn graph by computing unitigs from any set of strings in memory [23] or with external memory [33, 40, 41]. Incidentally, the set of unitigs computed from a set of strings is also a way to store a set of k -mers without repetition, and thus in reasonably small space. However, the necessity to preserve the topology of the graph makes unitigs an inferior choice to represent k -mer sets, as the sum of their length is still far from optimal, and they consist of many separate strings. The possibility to ignore the topology for k -mer-based methods opens more leeway in their representation that can be exploited to reduce the resource consumption of existing and future bioinformatics tools.

The need for such representations becomes apparent when observing the amount of data available to bioinformaticians. For example, the number of complete bacterial genomes available in RefSeq [42] more than doubled between May 2020 and July 2021 from around 9000 [43] to around 21,000. And with the ready availability of modern sequencing technologies, the amount of genomic data will increase further in the next years. In turn, analysing this data requires an ever growing amount of computational resources. But this could be relieved through a smaller representation that reduces the RAM usage and speeds up the analysis tools, and thereby allows to run larger pipelines using less computational resources. To fulfil this goal, a *plain text* representation would be the most useful: if the representation has to be decompressed before usage, then this likely erases the savings in RAM and/or adds additional runtime overhead. Formally, a plain text representation is a set of strings that contains each k -mer from the input strings (forward, reverse-complemented, or both) and no other k -mer. We denote such a set as a *spectrum preserving string set* (SPSS). Note that this definition is different from the SPSS defined by Rahman and Medvedev [44], who consider an SPSS to include the additional restriction that each k -mer must be present at most once. Such a plain text representation has the great advantage that some tools (like, e.g. Bifrost's query [23]) can use it without modification. We expect that even those tools that require modifications would not be hard to modify (like, e.g. SSHash [45] which we modified here as an example).

Related work

The concept of storing a set of k -mers in plain text without repeating k -mers to achieve a smaller and possibly simpler representation has recently been simultaneously discovered and named *spectrum preserving string sets [without k -mer repetition]* by Rahman and Medvedev [44] as well as *simplitigs* by Břinda, Baym and Kucherov [43]. To avoid confusion with our redefinition of the SPSS, we call this concept *simplitigs* in our work. Both Rahman and Medvedev and Břinda, Baym and Kucherov propose an implementation that greedily joins consecutive unitigs to compute such a representation. The UST algorithm by Rahman and Medvedev works on the node-centric de Bruijn graph of the input strings and finds arbitrary paths in the graph starting from arbitrary nodes. Each node is visited exactly by one path, and whenever a path cannot be extended forwards (because a dead-end was found, or all successor nodes have been visited already), then a new path is started from a new random node. Before a new path is started this way, if any successor node of the finished path marks the start of a different path, then the two

paths are joined. During the traversal, the unitigs of the visited nodes are concatenated (without repeating the $k - 1$ overlapping characters) and those strings are the final output. Břinda, Baym and Kucherov's greedy algorithm to compute simplitigs (for which the authors provide an implementation under the name ProphAsm [43]) does not construct a de Bruijn graph, but instead collects all k -mers into a hash table. Then it extends arbitrary k -mers both forwards and backwards arbitrarily until they cannot be extended anymore, without repeating any k -mers. The extended k -mers are the final output.

Both heuristics greatly reduce the number of strings (*string count*, SC) as well as the total amount of characters in the strings (*cumulative length*, CL) required to store a k -mer set. The reduction in CL directly relates to a lower memory consumption for storing a set of strings, but also the reduction in SC is very useful. For example, when storing a set of strings, not only the strings need to be stored, but also some index structure telling where they start and end. This structure can be smaller if less strings exist. There might even be cases where by increasing CL and decreasing SC, the overall size of the representation of the string set (strings + index structure) can be decreased. However, to stay independent of any specific data structure we only optimise CL. Břinda, Baym and Kucherov show that both SC and CL are greatly reduced for very tangled de Bruijn graphs, like graphs for single large genomes with small k -mer length and pangenome graphs with many genomes. Additionally they show merits of using heuristic simplitigs in downstream applications like an improvement in run time of k -mer queries using BWA [46], as well as a reduction in space required when storing heuristic simplitigs compressed with general-purpose compressors over storing unitigs compressed in the same way. Rahman and Medvedev show a significant reduction in SC and CL on various data sets as well, and also show a reduction in space required to store heuristic simplitigs over unitigs when compressed with general-purpose compressors. Khan et al. [40] also provide an overview of using heuristic simplitigs for various genomes, including also a human gut metagenome.

The authors of both papers also give a lower bound on the cumulative length of simplitigs, and show that their heuristics achieve representations with a cumulative length very close to the lower bound for typical values of k (31 for bacterial genomes and 51 for eukaryotic genomes). Břinda, Baym and Kucherov also experiment with lower values of k (< 20 for bacterial genomes and < 30 for eukaryotic genomes) which make the de Bruijn graph more dense to almost complete, and show that in these cases, their heuristic does not get as close to the lower bound as for larger values of k . Further, the authors of both papers consider whether computing minimum simplitigs without repeating k -mers might be NP-hard. This has recently been disproven by Schmidt and Alanko [47], and in fact simplitigs with minimum cumulative length can be computed in linear time by using a subset of the matchtig algorithm. Their algorithm constructs a bidirected arc-centric de Bruijn graph in linear time using a suffix tree, and then Eulerises it by inserting *breaking arcs*. It then computes a bidirected Eulerian circuit in the Eulerised graph and breaks it at all breaking arcs. The strings spelled by the resulting walks are the optimal simplitigs, named *Eulertigs*. Specifically, they leave out all parts from the matchtigs algorithm that relate to concatenating unitigs by repeating k -mers, and instead only concatenate consecutive unitigs in an optimal way. In line with the previous results about the lower bounds [43, 44], Eulertigs are only marginally smaller than the strings computed

by previous heuristics. All these suggest that no further progress is possible when k -mer repetitions are not allowed in a plain text representation.

There are already tools available that use simplitigs. The compacted de Bruijn graph builder cuttlefish2 [40] has an option to output simplitigs instead of maximal unitigs. A recent proposal for a standardised file format for k -mer sets explicitly supports simplitigs [48]. Also the k -mer dictionary SShash [45] uses simplitigs to achieve a smaller representation and to answer queries more efficiently. Here, the higher efficiency is achieved both by reducing the space required to store the k -mers themselves, but also due to the lower string count reducing the size of the index data structures on top. Further, a recent proposal to index genomic sequences as opposed to k -mer sets works with simplitigs without modification [49], and with minor extra book-keeping also for general SPSSs. In that work, the size of the SPSS is very minor compared to the size of the index, however, major components of the index may be smaller if the SPSS contains less strings, which can be achieved by using greedy matchtigs. Our algorithms were also integrated into the external-memory de Bruijn graph compactor GGCAT [41], which was easy to do (source: personal communication with Andrea Cracco).

For compressing multiple k -mer sets, in [50] the authors use an algorithm similar to ProphAsm and UST that separates the unique k -mer content of each set from the k -mer content shared with other sets. For compressing k -mer sets with abundances, the plain text representation REINDEER [51] uses substrings of unitigs with k -mers of similar abundance, called *monotigs*. In the wider field of finding small representations of k -mer sets that are not necessarily in plain text, there exists for example ESSCompress [52], which uses an extended DNA alphabet to encode similar k -mers in smaller space.

Our contribution

In this paper we propose the first algorithm to find an SPSS of *minimum* size (CL). Moreover, we show that a minimum SPSS with repeated k -mers is polynomially solvable, based on a many-to-many min-cost path query and a min-cost perfect matching approach. We further propose a faster and more memory-efficient greedy heuristic to compute a small SPSS that skips the optimal matching step, but still produces close to optimal results in CL, and even better results in terms of SC.

Our experiments over references and read datasets of large model organisms and bacterial pangenomes show that the CL decreases by up to 26% and the SC by up to 90% over UST or ProphAsm (on larger datasets, sometimes UST cannot be run because BCALM2 cannot be run, and sometimes ProphAsm cannot be run because it does not use external memory). Compared to unitigs, the CL decreases by up to 59% and SC by up to 97%. These improvements come often at just minor costs, as computing our small representation (which includes a run of BCALM2) takes less than twice as long than computing unitigs with BCALM2, and takes less than 37% longer in most cases. Even if the memory requirements for large read datasets increase, they stay within the limits of a modern server.

Finally we show that besides the smaller size of a minimum SPSS, it also has advantages in downstream applications. As an example of a k -mer-based method, we query our compressed representation with the tools SShash [45] and Bifrost [23]. These are

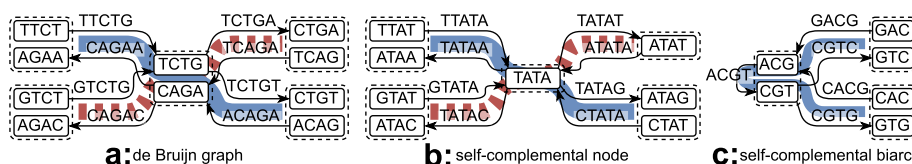


Fig. 1 Examples of de Bruijn graphs. Binodes are surrounded by a dashed box, where self-complemental binodes contain only one graph node. **a** A de Bruijn graph of the strings TTCTGA and GTCTGT. The colored and patterned lines are an arc-covering set of biwalks. **b** A de Bruijn graph containing two self-complemental nodes. The colored and patterned lines are a set of biwalks that visit each biarc exactly once. **c** A de Bruijn graph containing a self-complemental biarc. The bold colored line is a biwalk that visits each biarc exactly once

state-of-the-art tools supporting *k*-mer-based queries in genomic sequences, using a representation of a *k*-mer set as a set of unitigs. By simply replacing unitigs with the strings produced by our greedy heuristic, and without modifications to Bifrost and a minor modification to SSHash disabling features that require unique *k*-mers, we get a speedup of up to 4.26× over unitigs, and up to 2.10× over strings computed by UST and ProphAsm. We call the modified version of SSHash “SSHash-Lite”.

Results

Basic graph notation

We give detailed definitions for our notation below in the “Preliminaries” section, but give an intuition about the required notation for the results section here already. Note that our notation deviates from standard mathematical bidirected graph notations, but it is useful in practice as it allows to implement bidirected graphs on top of standard graph libraries. We assume that the reader is familiar with the general concept of de Bruijn graphs.

Our bidirected graphs are *arc-centric bidirected de Bruijn graphs*. Arc-centric de Bruijn graph means that *k*-mers are on the arcs, and nodes represent *k* – 1 overlaps between *k*-mers. We represent the bidirected graph as doubled graph, i.e. by having a separate forward and reverse arc for each *k*-mer and a separate forward and reverse node for each *k* – 1 overlap. In this graph, *binodes* are ordered pairs (v, v^{-1}) of nodes that are reverse complements of each other, and *biarcs* are ordered pairs (e, e^{-1}) of arcs that are reverse complements of each other. Two biarcs (e, e^{-1}) and (f, f^{-1}) are consecutive if the normal arcs *e* and *f* are consecutive, i.e. the *f* leaves the node entered by *e*. A *biwalk* is a sequence of consecutive biarcs. If a biwalk visits a biarc, then it is considered to be covering both directions of the biarc. See Fig. 1a for an example.

Matchtigs as a minimum plain text representation of *k*-mer sets

We introduce the *matchtig algorithm* that computes a character-minimum SPSS for a set of genomic sequences. While former heuristics (ProphAsm, UST) did not allow to repeat *k*-mers, our algorithm explicitly searches for all opportunities to reduce the character count in the SPSS by repeating *k*-mers. Consider for example the arc-centric de Bruijn graph in Fig. 2a. When representing its *k*-mers without repetition as in Fig. 2b, we need 43 characters and 7 strings. But if we allow to repeat *k*-mers as in Fig. 2d, we require only 39 characters and 5 strings. It turns out that structures similar to this

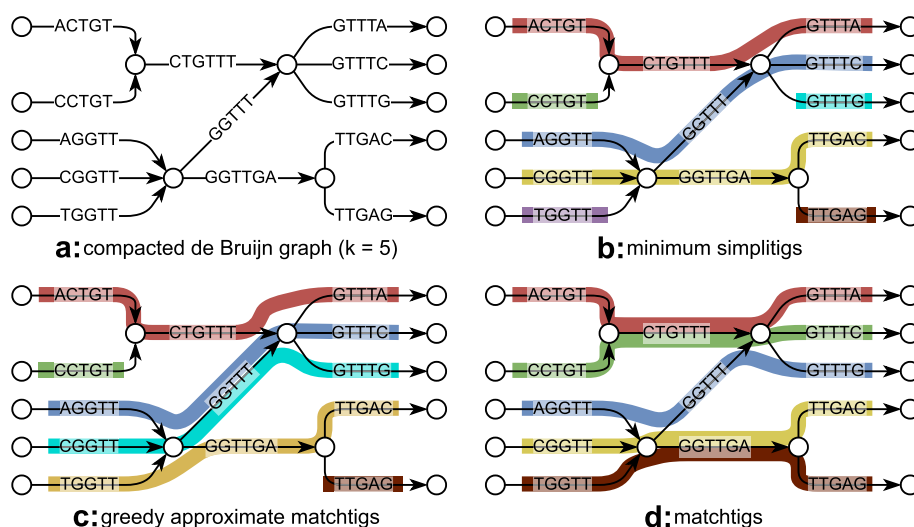


Fig. 2 Different SPSSs computed on an arc-centric de Bruijn graph $k = 5$. For simplicity, the reverse complements of all nodes and arcs are omitted. **a** The original de Bruijn graph in compacted form has 13 units with 70 total characters. **b** Example of simpltigs with 7 strings and 43 total characters. **c** Example of greedily approximated matchtigs with 6 strings and 40 total characters. **d** Example of matchtigs with 5 strings and 39 total characters

example occur often enough in real genome graphs to yield significant improvements in both character and string count of an SPSS.

Similar to previous heuristics, our algorithm works on the compacted bidirected de Bruijn graph of the input sequences. However, we require an arc-centric de Bruijn graph, but this can be easily constructed from the node-centric variant (see the “[Building a compacted bidirected arc-centric de Bruijn graph from a set of strings](#)” section). In this graph, we find a min-cost circular biwalk that visits each biarc at least once, and that can jump between arbitrary nodes at a cost of $k - 1$. This formulation is very similar to the classic Chinese postman problem [53], formulated as follows: find a min-cost circular walk in a directed graph that visits each arc at least once. This similarity allows us to adapt a classic algorithm from Edmonds and Johnson that solves the Chinese postman problem [54] (the same principle was applied in [55]). They first reduce the problem to finding a min-cost Eulerisation via a min-cost flow formulation, and then further reduce that to min-cost perfect matching using a many-to-many min-cost path query between unbalanced nodes. In a similar work [56], the authors solve the Chinese postman problem in a bidirected de Bruijn graph by finding a min-cost Eulerisation via a min-cost flow formulation. As opposed to [54, 55] and us, in [56] the authors propose to solve the min-cost flow problem directly with a min-cost flow solver. We believe this to be infeasible for our problem, since the arbitrary jumps between nodes require the graph in the flow formulation to have a number of arcs quadratic in the number of nodes.

Our resulting algorithm is polynomial but while it runs fast for large bacterial pangenomes, it proved practically infeasible to build the matching instance for very large genomes (≥ 500 Mbp). This is because each of the min-cost paths found translates into roughly one edge in the matching graph, and the number of min-cost paths raises quadratically if the graph gets denser. Thus, our algorithm ran out of memory when constructing it for larger genomes, and for those where we were able to construct the

matching instance, the matcher itself suffered from integer overflows, since it uses 32-bit integers to store the instance. Hence, for practical purposes, we introduce a greedy heuristic to compute approximate matchtigs. This heuristic does not build the complete instance of the matching problem, but just greedily chooses the shortest path from each unbalanced node to Eulerise the graph. This reduces the amount of paths per node to at most one, and as a result, the heuristic uses significantly less memory, runs much faster, and achieves near optimal speedups when run with multiple threads (see Additional file 1: Supplemental figure S1). While it can in theory produce suboptimal results as in Fig. 2c, in practice, the size of the greedily computed strings is very close to that of matchtigs, and the number of strings is always smaller.

Moreover, the minimality of matchtigs allows us to exactly compare, for the first time, how close heuristic algorithms to compute simplitigs are to optimal SPSS (on smaller genomes and on bacterial pangenomes, due to the resource-intensiveness of optimal matchtigs).

Our implementations are available on GitHub (<https://github.com/algbio/matchtigs>) as both a library and a command line tool, both written in Rust. They support both GFA and fasta file formats with special optimisations for fasta files produced by BCALM2 or GGCAT. Additionally, our implementations support gzip-compressed input and output, as well as outputting an ASCII-encoded bitvector of duplicate k -mers.

Compression of model organisms

We evaluate the performance of our proposed algorithms on three model organisms: *C. elegans*, *B. mori*, and *H. sapiens*. We benchmark the algorithms on both sets of short reads (average length 300 for *C. elegans* and *B. mori*, and 296 for *H. sapiens*) and reference genomes of these organisms. On human reads, we filter the data during processing so that we keep only k -mers that occur at least 10 times (min abundance = 10). This is because with a min abundance of 1, *H. sapiens* has 114 billion unique k -mers. This extreme k -mer count causes ProphAsm and our tool to run out of memory even with 2TiB of RAM.

We use the metrics cumulative length (CL) and string count (SC) as in [43]. The CL is the total number of characters in all strings in the SPSS, and the SC is the number of strings. We evaluate our algorithms against the same large genomes as in [43], using both the reference genome and a full set of short reads of the respective species (see Table 1 for the results). Since UST as well as matchtigs and greedy matchtigs require unitigs as input, and specifically UST needs some extra information in a format only output by BCALM2 [33], we run BCALM2 to compute unitigs from the input strings. We chose $k = 31$, as it is commonly used in k -mer-based methods. While for larger genomes, larger k are used as well, we use the value $k = 31$ throughout the main matter to allow for easier comparison between results. Furthermore, for all data sets but the *C. elegans* reference, the matchtigs algorithm ran out of the given 256GiB memory, so we only compute greedy matchtigs for those.

On read data sets where we keep all k -mers, our greedy heuristic achieves an improvement of up to 26% CL and 82% SC over the best competitor (UST-tigs). The human read data set has smaller improvements, however it was processed with a min abundance of

Table 1 Quality and performance of compressing model organisms

Genome	Algorithm	CL ratio	SC ratio	Time [s]	Memory [GiB]
<i>C. elegans</i> (reads)	B2	1.00	1.00	2402	5.54
	B2+UST	0.58	0.37	3424	(1.43) 17.6 (3.18)
	ProphAsm	0.55	0.34	5433	(2.26) 56.5 (10.2)
	B2+GREEDY	0.45 (0.79)	0.11 (0.28)	3057	(1.27) 41.0 (7.41)
<i>B. mori</i> (reads)	B2	1.00	1.00	6406	9.95
	B2+UST	0.55	0.35	9896	(1.54) 56.2 (5.64)
	ProphAsm	0.52	0.31	27,912	(4.36) 157 (15.8)
	B2+GREEDY	0.41 (0.74)	0.06 (0.18)	11,793	(1.84) 123 (12.4)
<i>H. sapiens</i> (reads)	B2	1.00	1.00	168,938	12.4
	B2+UST	0.67	0.46	170,427	(1.01) 29.0 (2.34)
	B2+GREEDY	0.57 (0.84)	0.22 (0.48)	209,646	(1.24) 68.5 (5.52)
<i>C. elegans</i>	B2	1.00	1.00	52.7	0.96
	B2+UST	0.92	0.34	58.6	(1.11) 0.96 (1.00)
	ProphAsm	0.92	0.30	133	(2.52) 3.78 (3.94)
	B2+GREEDY	0.90 (0.98)	0.06 (0.18)	59.9	(1.14) 0.96 (1.00)
	B2+MATCH	0.90 (0.98)	0.07 (0.23)	380	(7.21) 1.34 (1.40)
<i>B. mori</i>	B2	1.00	1.00	244	1.92
	B2+UST	0.78	0.34	303	(1.24) 1.92 (1.00)
	ProphAsm	0.76	0.28	716	(2.93) 13.8 (7.19)
	B2+GREEDY	0.72 (0.92)	0.06 (0.19)	334	(1.37) 2.42 (1.26)
<i>H. sapiens</i>	B2	1.00	1.00	1787	6.29
	B2+UST	0.79	0.33	2249	(1.26) 8.80 (1.40)
	ProphAsm	0.76	0.26	6677	(3.74) 130 (20.7)
	B2+GREEDY	0.71 (0.91)	0.03 (0.10)	4999	(2.80) 17.3 (2.75)

We chose $k = 31$ and a min abundance of 10 for *H. sapiens* reads and 1 for all others. The CL and SC ratios are between compressed strings and unitigs, and in parentheses are the ratios between our algorithm and the best competitor. B2 means BCALM2. For time and memory, we report the total time and maximum memory required to compute the tigs from the respective data set. BCALM2 directly computes unitigs and ProphAsm directly computes heuristic simpltigs. UST, GREEDY and MATCH compute heuristic simpltigs, greedy matchtigs and matchtigs from unitigs. The number in parentheses behind time and memory indicates the slowdown/increase over computing just unitigs with BCALM2. All algorithms were run with 28 threads, except for UST which supports only one thread (the preceding run of BCALM2 was still executed with 28 threads), and ProphAsm, which supports only one thread as well. Matchtigs are too expensive to run on all genomes except for the *C. elegans* reference, and ProphAsm takes too much time on *H. sapiens* reads, especially since it does not support minimum abundance. The lengths of the genomes are 100Mbp for *C. elegans*, 482Mbp for *B. mori* and 3.21Gbp for *H. sapiens* and the read data sets have a coverage of $64\times$ for *C. elegans*, $58\times$ for *B. mori* and $300\times$ for *H. sapiens*. The unique k -mer counts of the read datasets are 1.35 billion for *C. elegans*, 3.66 billion for *B. mori* and 2.78 billion for *H. sapiens*.

10, yielding longer unitigs with less potential for compression. On reference genomes, the improvement in CL is smaller with up to 7%; however, the improvement in SC is much larger with up to 90%.

For *C. elegans*, where computing matchtigs is feasible as well, we observe that they yield no significant improvement in CL, but are even slightly worse in SC than the greedy heuristic. The greedy heuristic actually optimises SC more than the optimal matchtigs algorithm. That is because the matching instance in the optimal algorithm is built to optimise CL, and whenever joining two strings does not alter CL, the choice is made arbitrarily. On the other hand, the greedy heuristic makes as many joins as possible, as long as a join does not worsen the CL. This way, the greedy heuristic actually prioritises joining two strings even if it does not alter the CL. For more details, see the “[Solving the min-cost integer flow formulation with min-cost matching](#)” and “[Efficient computation](#)”

of the greedy heuristic” sections. See Additional file 1: Supplemental figure S2 for more quality measurements with different k -mer size and min. abundance.

We assume that the improvements correlate inversely with the average length of maximal unitigs of the data set. Our approach achieves a smaller representation by joining unitigs with overlapping ends, avoiding the repetition of those characters. This has a natural limit of saving at most $k - 1$ characters per pair of unitigs joint together, so at most $k - 1$ characters per unitig. In turn, the maximum fraction of characters saved is bound by $k - 1$ divided by the average length of unitigs. In Additional file 1: Supplemental figure S2, we have varied the k -mer size and min. abundance for our data sets to vary the average length of unitigs. This gives us visual evidence for a correlation between average unitig length and decrease in CL.

Our improvements come at often negligible costs in terms of time and memory. Even for read sets, the run time at most doubles compared to BCALM2 in the worst case. However, the memory consumption rises significantly for read sets. This is due to the high number of unitigs in those graphs and the distance array of Dijkstra’s algorithm, whose size is linear in the number of nodes and the number of threads. See Additional file 1: Supplemental figure S3 for more performance measurements with different k -mer size and min. abundance.

Compression of pangenomes

In addition to model organisms with large genomes, we evaluate our algorithms on bacterial pangenomes of *N. gonorrhoeae*, *S. pneumoniae*, *E. coli*, and *Salmonella*, as well as a *human* pangenome. We use the same metrics as for model organisms. For the bacterial genomes, we choose $k = 31$, but also for the human genome for the reasons argued above, and also for easier comparability of the results on the different genomes. We show the results in Table 2. See Additional file 1: Supplemental figure S4 for more quality measurements with different k -mer size and min. abundance, and Additional file 1: Supplemental figure S5 for more performance measurements with different k -mer size and min. abundance. In neither of them, we have included *Salmonella* or human, as they take too much time.

Our algorithms improve CL up to 19% (using greedy matchtigs) over the best competitor and SC up to 70% (using greedy matchtigs). Matchtigs always achieve a slightly lower CL and slightly higher SC than greedy matchtigs, but the CL of greedy matchtigs is always at most 2% worse than that of matchtigs. We again assume that the improvements are correlated inversely to the average size of unitigs, as suggested by the experiments in Additional file 1: Supplemental figure S4. These improvements come at negligible costs, using at most 15% more time and 11% more memory than BCALM2 when computing greedy matchtigs, except for the large salmonella pangenome, which took 50% more memory. The higher memory consumption is due to the graph being more tangled due to the high number of genomes in the pangenome. For matchtigs, the time increases by less than a factor of three and memory by at most 12% compared to BCALM2. See Additional file 1: Supplemental figure S5 for more performance measurements with different k -mer size and min. abundance.

Table 2 Quality and performance of compressing pangenomes

Pangenome	Algorithm	CL ratio	SC ratio	Time [s]	Memory [GiB]
1102× <i>N. gonorrhoeae</i>	B2	1.00	1.00	29.1	4.25
	B2+UST	0.63	0.35	31.1 (1.07)	4.25 (1.00)
	ProphAsm	0.62	0.33	735 (25.3)	0.202 (0.05)
	B2+GREEDY	0.57 (0.93)	0.18 (0.54)	30.2 (1.04)	4.25 (1.00)
	B2+MATCH	0.57 (0.92)	0.18 (0.56)	31.1 (1.07)	4.25 (1.00)
616× <i>S. pneumoniae</i>	B2	1.00	1.00	26.1	3.07
	B2+UST	0.61	0.35	31.1 (1.19)	3.07 (1.00)
	ProphAsm	0.60	0.33	445 (17.0)	0.424 (0.14)
	B2+GREEDY	0.53 (0.89)	0.13 (0.41)	29.0 (1.11)	3.07 (1.00)
	B2+MATCH	0.52 (0.88)	0.14 (0.44)	41.8 (1.60)	3.07 (1.00)
3682× <i>E. coli</i>	B2	1.00	1.00	334	6.95
	B2+UST	0.60	0.35	417 (1.25)	6.95 (1.00)
	ProphAsm	0.59	0.32	13,339 (39.9)	7.05 (1.01)
	B2+GREEDY	0.51 (0.87)	0.11 (0.33)	384 (1.15)	6.95 (1.00)
	B2+MATCH	0.50 (0.85)	0.12 (0.37)	861 (2.58)	7.78 (1.12)
~309k× <i>Salmonella</i>	B2	1.00	1.00	82,417	12.7
	B2+UST	0.57	0.36	82,841 (1.01)	12.7 (1.00)
	B2+GREEDY	0.46 (0.81)	0.11 (0.30)	82,726 (1.00)	19.1 (1.50)
2505× <i>H. sapiens</i>	CF	1.00	1.00	77,582	402
	CF+ProphAsm	0.68	0.31	82,797 (1.07)	402 (1.00)
	CF+GREEDY	0.63 (0.93)	0.16 (0.50)	83,507 (1.08)	402 (1.00)

We chose $k = 31$ and a min abundance of 1. The CL and SC ratios are between compressed strings and unitigs, and in parentheses are the ratios between our algorithm and the best competitor. B2 means BCALM2. For time and memory, we report the total time and maximum memory required to compute the tigs from the respective data set. BCALM2 directly computes unitigs and ProphAsm directly computes heuristic simplitigs. UST, GREEDY and MATCH compute heuristic simplitigs greedy matchtigs and matchtigs from unitigs. The number in parentheses behind time and memory indicates the slowdown/increase over computing just unitigs with BCALM2. All algorithms were run with 28 threads, except for UST which supports only one thread (the preceding run of BCALM2 was still executed with 28 threads), and ProphAsm, which supports only one thread as well. The *N. gonorrhoeae* pangenome contains 8.36 million unique k -mers, the *S. pneumoniae* pangenome contains 19.3 million unique k -mers, the *E. coli* pangenome contains 341 million unique k -mers, the *Salmonella* pangenome contains 657 million unique k -mers and the human pangenome contains 2.8 billion unique k -mers. Due to its size, ProphAsm and MATCH could not be run on the *Salmonella* pangenome. Also due to size, BCALM2 did not run on the human pangenome, hence we used Cuttlefish 2. To still be able to compare against competitors, we ran ProphAsm on the unitigs produced by Cuttlefish 2 (UST requires extra information from BCALM2). To let Cuttlefish 2 run faster, we have used the flag `-unrestricted-memory`. Hence, its memory consumption is a lot higher than that of BCALM2

***k*-mer-based short read queries**

Matchtigs have further applications beyond merely reducing the size required to store a set of k -mers. Due to their smaller size and lower string count, they can make downstream applications more efficient. To make a concrete example, in this section we focus on *membership queries*. As already explained, each SPSS (unitigs, UST-tigs, matchtigs, etc.) can be considered as a (multi-) set of k -mers. Given a k -mer, a membership query is to verify whether the k -mer belongs to the set or not. We focus on *exact queries*, rather than approximate, i.e. if a k -mer does not belong to the set then the answer to the query *must* be “false”. Assessing the membership to the set for a string Q longer than k symbols is based on the answers to its constituent k -mers: only if *at least* $\lfloor \theta \times (|Q| - k + 1) \rfloor$ k -mers of Q belongs to the set, then Q is considered to be present in the set. The threshold θ is therefore an “inclusion” rate, which we fix to 0.8 for the experiments in this section.

To support fast membership queries in compressed space, we build an SSHash-Lite dictionary over each SPSS. SSHash-Lite is a relaxation of SSHash [45, 57] in that it supports membership queries *without* requiring each k -mer to appear once in the underlying SPSS. It is available at <https://github.com/jermp/sshash-lite>. In short, SSHash is a compressed dictionary for k -mers — based on minimal perfect hashing [58] and minimizers [59] — which, for an input SPSS without duplicates and having n (distinct) k -mers, assigns to each k -mer in the input a unique integer number from 0 to $n - 1$ by means of a Lookup query. The result of Lookup for any k -mer that is *not* contained in the input SPSS is -1 . Therefore, SSHash serves the same purpose of a minimal perfect hash function over a SPSS but it is also able to reject alien k -mers. Two variants of SSHash were proposed — a *regular* and a *canonical* one. The canonical variant uses some extra space compared to the regular one but queries are faster to answer. (For all further details, we point the reader to the original papers [45, 57].)

Now, to let SSHash be able to query SPSSs with possible duplicate k -mers (e.g. matchtigs), it was only necessary to modify the return value of the Lookup query to just return “true” if a k -mer is found in the dictionary rather than its unique integer identifier (respectively, “false” if a k -mer is not found instead of -1). Therefore, SSHash-Lite can be directly used to index and query the unitigs, UST-tigs, and matchtigs as well.

We compare the performance of SSHash-Lite when indexing unitigs, UST-tigs, and matchtigs in Table 3. We build the SPSSs from three datasets: a $\sim 309\text{k} \times$ *Salmonella* Enterica pangenome; a $300 \times$ coverage human short read dataset filtered to exclude k -mers with an abundance lower than 10; and a $2505 \times$ human pangenome. The *Salmonella* pangenome was queried with 3 million random *Salmonella* short reads with lengths between 70 and 502, and an N75 of 302. The human queries for both the human read dataset and the human pangenome are 3 million random short reads (296 bases each) from the human read dataset.

We see that matchtigs improve the performance of membership queries in *both space and time* compared to unitigs and UST-tigs. While the difference is more evident when compared to unitigs, matchtigs also consistently outperform UST-tigs — achieving the lowest space usage and faster query time across almost all combinations of dataset and index variant (regular/canonical).

Note again that the speed up in searching time is more evident on the human reads dataset since it is much larger than the *Salmonella* pan-genome and it is generally less evident for the canonical index variant of SSHash-Lite because it is approximately $2 \times$ faster to query than the regular one. Remarkably, regular SSHash-Lite over matchtigs achieves 26 – 59% reduction in space over unitigs while being also $4.26 \times$ faster to query on the human reads datasets. Compared to UST-tigs instead, matchtigs still retain $2.10 \times$ faster query time while improving space by up to 8%. These results were achieved on a typical bioinformatics compute node with many logical cores (256) and a large amount of RAM (2TB). In Additional file 1: Supplemental table S2, we performed the same experiment on a server with focus on single-thread performance, achieving slightly smaller improvements.

The reduction in index space when indexing matchtigs is to be attributed to the lower string count and fewer nucleotides in the collection. The speedups achieved by SSHash-Lite when indexing matchtigs instead of unitigs can be explained as follows. When

Table 3 Performance characteristics of querying different tigs with SSHash-Lite

Genome	Algorithm	Index time	Search time	Search speedup	Index size	Size imprv.
		[min]	[sec]		[GiB]	
(a) regular SSHash-Lite						
~309k× <i>Salmonella</i> (0.75)	unitigs	2.77	3027	1.00	1.04	1.00
	UST	2.42	1491	2.03	0.71	1.48
	gMatchtigs	2.60	710	4.26 (2.10)	0.65	1.60 (1.08)
Human reads (0.75)	unitigs	18.9	558	1.00	4.60	1.00
	UST	17.1	499	1.12	3.63	1.27
	gMatchtigs	19.2	384	1.45 (1.30)	3.47	1.33 (1.05)
2505× Human (0.65)	unitigs	15.1	515	1.00	3.63	1.00
	ProphAsm	14.0	421	1.22	2.86	1.27
	gMatchtigs	14.8	363	1.42 (1.16)	2.86	1.27 (1.00)
(b) canonical SSHash-Lite						
~309k× <i>Salmonella</i> (0.75)	unitigs	3.94	1576	1.00	1.13	1.00
	UST	3.30	961	1.64	0.78	1.44
	gMatchtigs	3.71	572	2.75 (1.68)	0.74	1.52 (1.06)
Human reads (0.75)	unitigs	25.0	373	1.00	5.02	1.00
	UST	23.4	324	1.15	4.05	1.24
	gMatchtigs	26.3	266	1.40 (1.22)	3.94	1.28 (1.03)
2505× Human (0.65)	unitigs	21.3	340	1.00	4.26	1.00
	ProphAsm	20.1	258	1.32	3.48	1.22
	gMatchtigs	21.1	232	1.46 (1.11)	3.52	1.21 (0.99)

SSHash-Lite is run with $k = 31$ and a k -mer-inclusion rate of 0.8. On the *Salmonella* pan-genome, we used a minimizer length of 17 for the regular index and a minimizer length of 16 for the canonical index. On the human reads, we used a minimizer length of 20 for the regular index and a minimizer length of 19 for the canonical index. On the human pangenome, we used a minimizer length of 19 for the regular index and a minimizer length of 20 for the canonical index. The search speedup is with respect to unitigs, and the search speedup in parentheses is with respect to the strings computed by UST. Index time is the end-to-end time required to build the SSHash-Lite index: it includes reading the collections from disk and building the data structure using external memory. Searching time is the time required to check which reads have at least 80% of their k -mers in the input SPSS. The number in parentheses under the genome is the k -mer hitrate, i.e. the fraction of k -mers from the query that are part of the queried dataset

querying, SSHash-Lite streams through the k -mers of the query. At the beginning, the tig containing the first k -mer of the query is determined using a minimal perfect hash function over the minimizers of the input SPSS, as well as the position of the k -mer in the tig. For the subsequent k -mers of the query, SSHash-Lite attempts to “extend” the matching of the k -mer against the identified tig by just comparing the single nucleotide following the previous k -mer in the tig. Extending a match in this way is extremely fast not only because just a single nucleotide needs to be compared but also because it is a very cache-friendly algorithm, dispensing random accesses to the index entirely. However, each time an extension is not possible (either because we have a mismatch or we have reached the end of the current tig) a “full” new search is made in the index. The search consists in evaluating the minimal perfect hash function and locating the k -mer inside another tig. Clearly, a search is much more expensive due to cache misses compared to an extension. Now, using longer tigs with a lower tig count — the case for the matchtigs — increases the chance of extension, or equivalently, decreases the number of full searches in the index. Compared to UST-tig, matchtigs can be faster to query exactly because allowing repeated k -mers to appear in the tigs further helps in creating

opportunities for extension. Therefore, by reducing the number of full searches, we can reduce the overall runtime of the query.

See Additional file 1: Section 3 for a similar query experiment with the slightly older tool Bifrost.

Discussion

k -mer-based methods have found wide-spread use in many areas of bioinformatics over the past years. However, they usually rely on unitigs to represent the k -mer sets, since they can be computed efficiently with standard tools [23, 33, 40, 41]. Unitigs have the additional property that the de Bruijn graph topology can easily be reconstructed from them, since they do not contain branching nodes other than on their first and last k -mer. However, this property is not usually required by k -mer-based methods, which has opened the question if a smaller set of strings other than unitigs can be used to represent the k -mer sets. If such a representation was in plain text, it should be usable in most k -mer-based tools, by simply feeding it to the tool instead of unitigs.

Previous work has relaxed the unitig requirement of the representation of the k -mer sets to arbitrary strings without k -mer repetitions. This resulted in a smaller representation, leading to improvements in downstream applications. Additionally, previous work considered whether that finding an optimal representation without repeated k -mers is NP-hard, which was then disproven and shown to be linear-time solvable. We have shown that by allowing repetitions, there is a polynomial optimal algorithm that achieves better compression and improvements in downstream applications.

Conclusions

Our *optimum* algorithm compresses the representation significantly more than previous work. For practical purposes, we also propose a greedy heuristic that achieves near-optimum results, while being suitable for practical purposes in runtime and memory. Specifically, our algorithms achieve a decrease of 26% in size and 90% in string count over UST. Additionally, we have shown that our greedy representation speeds up downstream applications, giving an example with a factor of 2.10 compared to previous compressed representations.

Our implementation is available as a stand-alone command-line tool and as a library. We hope that our efficient algorithms result in a wide-spread adoption of near-minimum plain-text representations of k -mer sets in k -mer-based methods, resulting in more efficient bioinformatics tools.

Methods

We first give some preliminary definitions in the “[Preliminaries](#)” section and define our problem in the “[Problem overview](#)” section. Note that to stay closer to our implementation, our definitions of bidirected de Bruijn graphs differ from those in e.g. [56]. However, the concepts are fundamentally the same. Then, in the “[Building a compacted bidirected arc-centric de Bruijn graph from a set of strings](#)”, “[Reduction to the bidirected partial-coverage Chinese postman problem](#)”, “[Solving the bidirected partial-coverage Chinese postman problem with min-cost integer flows](#)”, “[Solving the min-cost integer flow formulation with min-cost matching](#)”, and “[Efficient computation of many-to-many](#)”

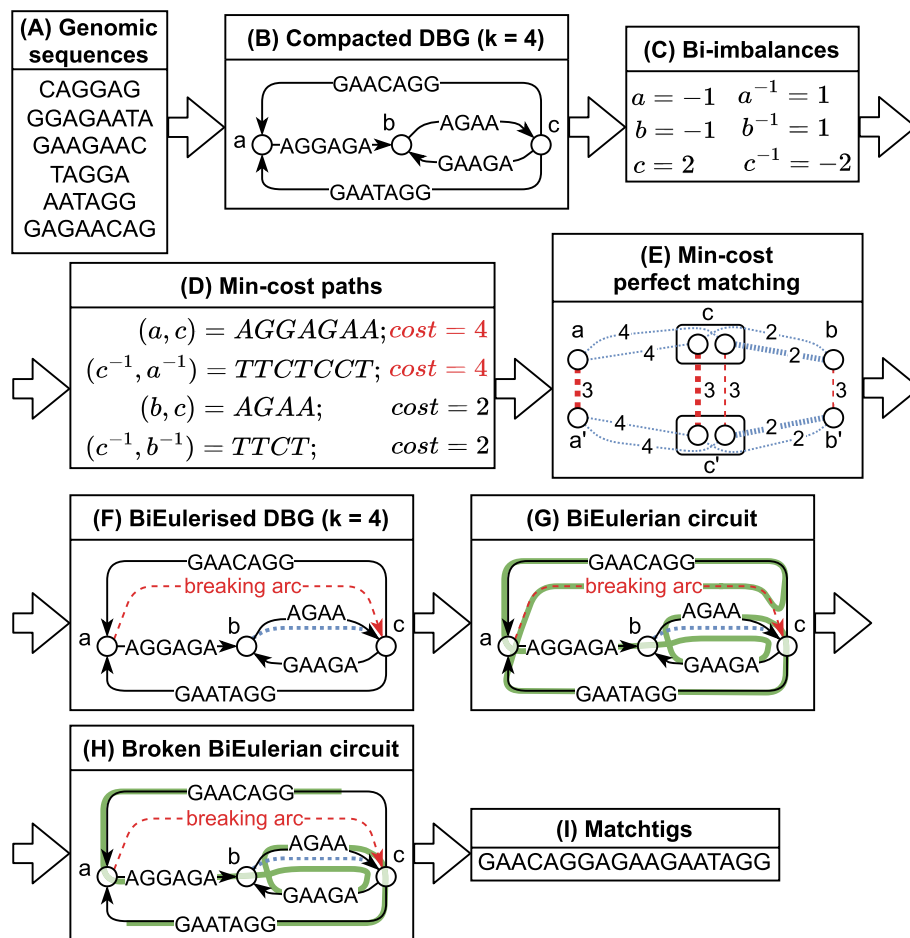


Fig. 3 An example of the matchtigs algorithm. **A** The input genomic sequences. **B** We first build an arc-centric compacted de Bruijn graph (for simplicity, the reverse complements of the nodes and arcs are not shown). **C** In the graph we compute the bi-imbances of the nodes (the difference between outdegree and indegree). **D** From each node with negative bi-imbalance we compute the min-cost paths to all reachable nodes with positive bi-imbalance. The costs of each arc are the amount of characters required to join two strings from the negative to the positive node while repeating the k -mers between the nodes. Specifically, the costs of an arc are $|s| - (k - 1)$, where $|s|$ is the length of its label. **E** Using a min-cost perfect matching instance built from the min-cost paths, we decide which bi-imbances should be fixed by repeating k -mers. The blue/tightly dashed edges are joining edges stemming from the min-cost paths. The red edges in longer dashes indicate that a node should stay unmatched, i.e. that fixing its bi-imbalance requires breaking arcs. The solution edges are highlighted in bold. There is one node in the matching problem for each binode in the original graph. The nodes x' are not reverse complements of nodes x , but stem from a reduction that makes a copy of each node. For more details, refer to the “Solving the min-cost integer flow formulation with min-cost matching” section. **F** For each joining edge in the solution we insert a joining arc into the DBG (in blue, small dashes), always directed such that the overall bi-imbalance decreases. The remaining imbalance is removed by inserting arbitrary breaking arcs (in red, longer dashes). **G** We compute a biEulerian circuit in the balanced graph. **H** We break the biEulerian circuit at all breaking arcs. **I** We output the strings spelled by the broken walks

“min-cost paths” sections, we describe how to compute matchtigs. The whole algorithm is summarised by an example in Fig. 3. For simplicity, we describe the algorithm using an uncompact de Bruijn graph. However, in practice it is much more efficient to use a compacted de Bruijn graph, but our algorithm can be adapted easily: simply replace the

costs of 1 for each original arc with the number of uncompact arcs it represents. In the “Efficient computation of the greedy heuristic” section, we describe the greedy heuristic.

Preliminaries

We are given an alphabet Γ and all strings in this work have only characters in Γ . Further, we are given a bijection $\text{comp} : \Gamma \rightarrow \Gamma$. The *reverse complement* of a string s is $s^{-1} := \text{rev}(\text{comp}^*(S))$ where rev denotes the reversal of a string and comp^* the character-wise application of comp . For an integer k , a string of length k is called a *k-mer*. From here on, we only consider strings of lengths at least k , i.e. strings that have at least one k -mer as substring. We denote the prefix of length $k - 1$ of a k -mer s by $\text{pre}(s)$ and its suffix of length $k - 1$ by $\text{suf}(s)$. The *spectrum* of a set of strings S is defined as the set of all k -mers and their reverse complements that occur in at least one string $s \in S$, formally $\text{spec}_k(S) := \{r \in \Gamma^k \mid \exists s \in S : r \text{ or } r^{-1} \text{ is substring of } s\}$.

An *arc-centric de-Bruijn graph* (or short *de-Bruijn graph*) $\text{DBG}_k(S) = (V, E)$ of order k of a set of strings S is defined as a standard directed graph with nodes $V := \{s \mid s \in \text{spec}_{k-1}(S)\}$ and arcs $E := \{(\text{pre}(s), \text{suf}(s)) \mid s \in \text{spec}_k(S)\}$. On top of this, we use the following notions of bidirectedness. An ordered pair of reverse-complementary nodes $[v, v^{-1}] \in V \times V$ is called a *binode* and an ordered pair of reverse-complementary arcs $[(a, b), (b^{-1}, a^{-1})] \in E \times E$ is called a *biarc*. Even though these pairs are ordered, reversing the order still represents the same binode/biarc, just in the other direction. A node v is called *canonical* if v is lexicographically smaller than v^{-1} , and an arc (a, b) is called *canonical* if the k -mer corresponding to (a, b) is lexicographically smaller or equal to the k -mer corresponding to (b^{-1}, a^{-1}) . If an arc or a node is its own reverse-complement (called *self-complemental*), then it is written as biarc $[(a, b)]$ or binode $[v]$. See Fig. 1 for examples of different bigraphs.

Since de Bruijn graphs are defined as standard directed graphs, we use the following standard definitions. The set of incoming (outgoing) arcs of a node is denoted by $E^-(v)$ ($E^+(v)$), and the indegree (outdegree) is $d^-(v) := |E^-(v)|$ ($d^+(v) := |E^+(v)|$). A *walk* in a de Bruijn graph is a sequence of adjacent arcs (followed in the forward direction) and a *unitig* is a walk in which all inner nodes (nodes with at least two incident walk-arcs) have exactly one incoming and one outgoing arc. The length $|w|$ of a walk w is the length of the sequence of its arcs (counting repeated arcs as often as they are repeated). A *compact de-Bruijn graph* is a de Bruijn graph in which all maximal unitigs have been replaced by a single arc. A *circular walk* is a walk that starts and ends in the same node, and a *Eulerian circuit* is a circular walk that contains each arc exactly once. A graph that admits a Eulerian circuit is *Eulerian*.

Assuming the complemental pairing of nodes and arcs defined above, we can define the following bidirected notions of walks and standard de Bruijn graph concepts. *Biwalks* and *circular biwalks* are defined equivalently to walks, except that they are sequences of biarcs. A biwalk w in a de Bruijn graph spells a string $\text{spell}(w)$ of overlapping visited k -mers. That is, $\text{spell}(w)$ is constructed by concatenating the string a from w 's first biarc $[(a, b), (b^{-1}, a^{-1})]$ (or $[(a, b)]$) with the last character of b of the first and all following biarcs. See Fig. 1 for examples of bidirected de Bruijn graphs and notable special cases.

A bidirected graph is *connected*, if between each pair of distinct binodes $[u, u^{-1}]$ and $[v, v^{-1}]$ that are not reverse complements of each other, there is a biwalk from $[u, u^{-1}]$ to $[v, v^{-1}]$ or from $[u, u^{-1}]$ to $[v^{-1}, v]$. We assume that our graph is connected, as on multiple disconnected components, our algorithm can be executed on each component, yielding a minimum result.

Problem overview

We are given a set of input strings I where each string has length at least k , and we want to compute a minimum spectrum preserving string set, defined as follows.

Definition 1

A *spectrum preserving string set* (or *SPSS*) of I is a set S of strings of length at least k such that $\text{spec}_k(I) = \text{spec}_k(S)$, i.e. both sets of strings contain the same k -mers, either directly or as reverse complement.

Note that our definition allows k -mers and their reverse complements to be repeated in the SPSS, both in the same string and in different strings. This is an important difference to the definition of an SPSS by Rahman and Medvedev [44]. Their definition is equivalent to *simplitigs*, defined by Břinda, Baym and Kucherov [43]; hence, we use the term *simplitigs* when we refer to an SPSS without k -mer repetitions in this paper. *Simplitigs* are an SPSS such that for each present k -mer the reverse complement is only present if the k -mer is self-complemental, and the k -mer is only present once in at most one string of the SPSS.

Definition 2

The *size* $||S||$ of an SPSS S is defined as

$$||S|| = \sum_{s \in S} |s|,$$

where $|s|$ denotes the length of string s . A *minimum* SPSS is an SPSS of minimum size.

On a high level, our algorithm works as follows (see also Fig. 3).

- 1 Create a bidirected de Bruijn graph from the input strings (see the “[Building a compacted bidirected arc-centric de Bruijn graph from a set of strings](#)” section).
- 2 Compute the *bi-imbances* of each node (see the “[Solving the bidirected partial-coverage Chinese postman problem with min-cost integer flows](#)” section).
- 3 Compute the min-cost bipaths of length at most $k - 1$ from all nodes with negative bi-imbalance to all nodes with positive bi-imbalance (see the “[Efficient computation of many-to-many min-cost paths](#)” section).
- 4 Solve a min-cost matching instance with costs for unmatched nodes to choose a set of shortest bipaths with minimum cost by reduction to min-cost perfect matching (see the “[Solving the min-cost integer flow formulation with min-cost matching](#)” section).
- 5 *BiEulerise* the graph with the set of bipaths as well as arbitrary arcs between unmatched nodes.

- 6 Compute a *biEulerian circuit* in the resulting graph (see the “[Solving the bidirected partial-coverage Chinese postman problem with min-cost integer flows](#)” section).
- 7 Break the circuit into a set of biwalks and translate them into a set of strings, which is the output minimum SPSS (see the “[Reduction to the bidirected partial-coverage Chinese postman problem](#)” section).

Note that in our implementation, a substantial difference is that we do not build the de Bruijn graph ourselves, but we expect the input to be a de Bruijn graph already. For our experiments, we use a compacted de Bruijn graph computed with BCALM2. We motivate the reasons for optimality while explaining our algorithm, but also give a formal proof in Additional file 1: Section 4.

Building a compacted bidirected arc-centric de Bruijn graph from a set of strings

When building the graph we first compute unitigs from the input strings using BCALM2. Then we initialise an empty graph and do the following for each unitig:

- 1 We insert the unitig’s first $k - 1$ -mer and its reverse complement as binode by inserting the two nodes separately and marking them as a bidirected pair, if it does not already exist. The existence is tracked with a hashmap, storing the two nodes corresponding to a k -mer and its reverse complement if it exists.
- 2 We do the same for the last $k - 1$ -mer of the unitig.
- 3 We add a biarc between the two binodes by inserting one forward arc between the forward nodes of the binodes, and one reverse arc between the reverse complement nodes of the binodes. The forward arc is labelled with the unitig, and the reverse arc is labelled with its reverse complement.

To save memory, we store the unitigs in a single large array, where each character is encoded as two-bit number. The keys of the hashmap and the labels of the arcs are pointers into the array, together with a flag for the reverse complement. Nodes do not need a label, as their label can be inferred from any of its incident arcs’ label. Recall that in the description of our algorithm, we use an uncompact graph only for simplicity.

Reduction to the bidirected partial-coverage Chinese postman problem

We first compute the arc-centric de-Bruijn graph $\text{DBG}_k(I)$ of the given input string set I as described in the “[Building a compacted bidirected arc-centric de Bruijn graph from a set of strings](#)” section. In $\text{DBG}_k(I)$, an SPSS S is represented by a *biarc-covering* set of biwalks W (the reverse direction of a biarc does not need to be covered separately). That is a set of biwalks such that each biarc is element of at least one biwalk (see Fig. 1a). According to the definition of spell, the size of S is related to W as follows:

$$||S|| = \sum_{w \in W} |\text{spell}(w)| = |W|(k - 1) + \sum_{w \in W} |w|. \quad (1)$$

Each walk costs $k - 1$ characters because it contains the node a from its first biarc $[(a, b), (b^{-1}, a^{-1})]$ (or $[(a, b)]$), and it additionally costs one character per arc.

To minimise $||S||$, we transform the graph as follows:

Definition 3

(Graph transformation) Given an arc-centric de-Bruijn graph $DBG_k(I) = (V, E)$, the *transformed graph* is defined as $DBG'_k(I) = (V, E')$ where E' is a multiset defined as $E' := E \cup (V \times V)$. In E' , arcs from E are marked as *non-breaking*, and arcs from $V \times V$ are marked as *breaking* arcs. The *cost function* $c(e), e \in E'$ assigns all non-breaking arcs the costs 1 and all breaking arcs the costs $k - 1$.

The reverse complements of breaking arcs are breaking as well, and the same holds for non-breaking arcs. This means that biarcs always are either a pair of reverse complementary breaking arcs, in which case we call them *breaking biarcs*, or a pair of reverse complementary non-breaking arcs, in which case we call them *non-breaking biarcs*. By the same pattern, self-complementary biarcs are defined to be *breaking biarcs* or *non-breaking biarcs* depending on their underlying arc. Breaking arcs have the costs $k - 1$ because each breaking arc corresponds to starting a new walk, which by Eq. 1 costs $k - 1$.

In the transformed graph we find a circular biwalk w^* of minimum cost that covers at least all original biarcs (to cover a biarc it is enough to traverse it once in one of its directions), as well as at least one breaking biarc. The reason for having at least one breaking biarc is that later we break the circular original-biarc-covering biwalk at all breaking biarcs to get a set of strings. If there was no breaking biarc, then we would get a circular string. Simply breaking a circular string at an arbitrary binode or a repeated non-breaking biarc does not produce a minimum solution in general, because there may be multiple circular original-biarc-covering biwalks with minimum costs, but with different repeated k -mers. When breaking the walk by removing a longer sequence of repeated k -mers, the resulting string gets shorter, the more repeated k -mers get removed. Hence we make finding an optimal breaking point part of the optimisation problem. We define such a walk as:

Definition 4

Given a transformed graph $DBG'_k(I) = (V, E)$, a *circular original-biarc-covering biwalk* is a circular biwalk w such that for each non-breaking arc $(a, b) \in E$ there is a biarc $[(a, b), (b^{-1}, a^{-1})]$, $[(b^{-1}, a^{-1}), (a, b)]$ or $[(a, b)]$ in w . Additionally, w needs to contain at least one breaking biarc.

Definition 5

Given a transformed graph $DBG'_k(I)$ and a circular original-biarc-covering walk w possibly consisting of biarcs $[(a, b), (b^{-1}, a^{-1})]$ and self-complementary biarcs $[(a, b)]$. The *costs* $c(w)$ of w are the sum of the costs of each biarc and self-complementary biarc, where the costs of a biarc $[(a, b), (b^{-1}, a^{-1})]$ are $c((a, b))$, and the costs of a self-complementary biarc $[(a, b)]$ are $c((a, b))$.

This is similar to the directed Chinese postman problem (DCPP). In the DCPP, the task is to find a circular min-cost arc-covering walk in a directed graph. It is a classical problem, known to be solvable in $O(n^3)$ time [60] with a flow-based algorithm, using e.g. [61] to compute min-cost flows. The partial coverage variant of the DCPP (requiring to cover only a subset of the arcs) is also known as the rural postman problem [62]. Further, the bidirected variant of the DCPP was discussed before in [56], and the authors also solved it using min-cost flow in $O(n^2 \log^2(n))$ time.

We break the resulting min-cost circular original-biarc-covering biwalk w^* at all breaking arcs (hence we require it to contain at least one breaking biarc, otherwise we would get a circular string). The returned SPSS is minimum, because the metric optimised when finding w^* matches Eq. 1.

Solving the bidirected partial-coverage Chinese postman problem with min-cost integer flows

Edmonds and Johnson [54] introduced a polynomial-time flow-based approach that is adaptable to solve our variant of the DCP. They show that finding a minimum circular arc-covering walk in a directed graph is equivalent to finding a minimum *Eulerisation* of the graph, and then any Eulerian circuit in the Eulerised graph. A Eulerisation is a multiset of arc-copies from the graph that makes the graph Eulerian if added, either by connecting nodes with missing outgoing arcs directly to nodes with missing incoming arcs, or by connecting them via a path of multiple arcs. A minimum Eulerisation is one whose sum of arc costs is minimum among all such multisets. To find such a minimum cost set of arcs, they formulate a min-cost flow problem as an integer linear program as follows:

$$\min \sum_{e \in E} c_e x_e \text{ s.t.} \quad (2)$$

$$x_e \text{ are non-negative integers, and} \quad (3)$$

$$\forall v \in V : \sum_{e \in E^-(v)} x_e - \sum_{e \in E^+(v)} x_e = d^+(v) - d^-(v). \quad (4)$$

The variable x_e is interpreted as the amount of flow through arc e , and the variable c_e denotes the costs for assigning flow to an arc e . The costs are equivalent to the arc costs in the weighted graph specified by the DCP instance. Objective (2) minimises the costs of inserted arcs as required. To ensure that the resulting flow can be directly translated into added arcs, Condition (3) ensures that the resulting flow is non-negative and integral. Lastly, Eq. (4) is the *balance constraint*, ensuring that the resulting flow is a valid Eulerisation of the graph. This constraint makes nodes with missing outgoing arcs into *sources*, and nodes with missing incoming arcs into *sinks*, with demands matching the number of missing arcs. Note that in contrast to classic flow problems, this formulation contains no capacity constraint. For a solution of this linear program, the corresponding Eulerisation contains x_e copies of each arc e .

To adapt this formulation to our variant of the DCP, we need to make modifications, namely:

- compute the costs while treating self-complemental biarcs the same as other biarcs,
- allow for partial coverage,
- force cover at least one breaking arc (one of the arcs that are not required to be covered),
- adjust the balance constraint for biwalks and
- ensure that the resulting flow is *bidirected*, i.e. the flow of each arc equals the flow of its reverse complement.

Bidirected costs

If we would simply count the costs of each arc separately, then self-complemental biarcs would cost 1 for each repetition, while other biarcs would cost 2 for each repetition, since other biarcs consist for two arcs. To circumvent this, we only count arcs corresponding to canonical k -mers in the cost function:

$$\min \sum_{e \in E, e \text{ is canonical}} c_e x_e \text{ s.t.} \quad (5)$$

Partial coverage

In the partial coverage Chinese postman problem, we are additionally given a set $F \subseteq E$ of arcs to be covered. In contrast to the standard DCP, a solution walk only needs to cover all the arcs in F . In our case, the set F is the set of original arcs of the graph before Eulerisation. To solve the partial coverage Chinese postman problem we define outgoing covered arcs $F^+(v) := F \cap E^+(v)$, and incoming covered arcs $F^-(v) := F \cap E^-(v)$ for a node v , as well as the covered outdegree $d_F^+(v) := |F^+(v)|$ and the covered indegree $d_F^-(v) := |F^-(v)|$. Then we reformulate the balance constraint as:

$$\forall v \in V : \sum_{e \in E^-(v)} x_e - \sum_{e \in E^+(v)} x_e = d_F^+(v) - d_F^-(v).$$

The resulting set of arcs is a minimum Eulerisation of the graph (V, F) , and a Eulerian walk in this graph is equivalent to a minimum circular F -covering walk in the original graph.

Cover one breaking arc

Next to the partial coverage condition, we additionally require to cover at least one of the arcs that is not required to be covered. Since we forbid negative flows, we can express this as:

$$\sum_{e \in (E \setminus F)} x_e \geq 1. \quad (6)$$

Bidirected balance

In contrast to Edmonds and Johnson, we are interested in a minimum circular *biwalk* that covers each original *biarc*. Analogue to the formulation for directed graphs, we make the following definitions:

Definition 6

(BiEulerian circuits and graphs) A *biEulerian circuit* in a bidirected graph is a bidirected circuit that visits each biarc exactly once. A *biEulerian graph* is a graph that admits a biEulerian circuit. A *biEulerisation* is a multiset of biarc-copies from the graph that makes a graph biEulerian if added. A *minimum biEulerisation* is one whose sum of arc costs is minimum among all biEulerisations of the same graph.

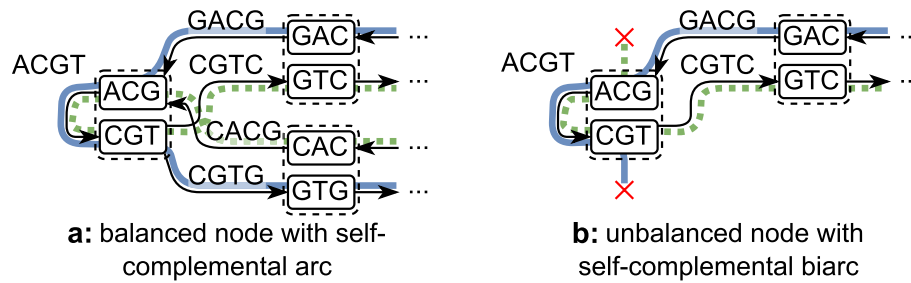


Fig. 4 Examples for self-complemental nodes and arcs. A self-complemental biarc $[(ACG, CGT)]$ covered by a biEulerian circuit $([(GAC, ACG), (CGT, GTC)], [(ACG, CGT)], [(CGT, GTG), (CAC, ACG)])$. The two directions of the bidirected circuit are drawn in blue $([(GAC, ACG), (ACG, CGT), (CGT, GTG)])$ and green dotted $([(CAC, ACG), (ACG, CGT), (CGT, GTC)])$. In **a**, the binode $[ACG, CGT]$ is balanced, hence the circuit can enter it with some biarc, cover the self-complemental biarc, and then leave it via some other biarc. If, like in **b**, there was no other biarc to leave $[ACG, CGT]$, then the graph would not be biEulerian, as the biarc $[(GAC, ACG), (CGT, GTC)]$ cannot be used twice, even if the second use is in the other direction. Visually, the blue walk cannot use (CGT, GTC) , since it was already used by the green dotted walk, and the green dotted walk cannot use (GAC, ACG) , as it was already used by the blue walk

We can compute a biEulerisation in the same way as we compute a Eulerisation, the only change is in the balance constraint. Observe that for a Eulerian graph, the *imbalance* $i_v := d^-(v) - d^+(v)$ is zero for each node [63], because each node is entered exactly as often as it is exited. For binodes, the definition of the *bi-imbalance* bi_v of a binode $[v, v^{-1}]$ or $[v]$ follows the same idea. However, in contrast to directed graphs, there are two special cases. These are mutually exclusive, since the labels of an arc are of length k and those of a node are of length $k - 1$, such that only one can have even length, which is required to be self-complemental in DNA alphabets.

Binodes $[v, v^{-1}] \in V \times V$ with $v \neq v^{-1}$ may have incident self-complemental arcs $[(v, v^{-1})]$ and/or $[(v^{-1}, v)]$ (see Fig. 1c for an example). If e.g. only $[(v, v^{-1})]$ exists, then to traverse it, a biwalk needs to enter v twice. First, it needs to reach $[v, v^{-1}]$ via some biarc, and after traversing $[(v, v^{-1})]$, it needs to leave $[v^{-1}, v]$ via a different biarc, whose reverse complement enters $[v, v^{-1}]$. Hence, a self-complemental biarc alters the bi-imbalance of a node by two. See Fig. 4 for an example of this situation. If only $[(v^{-1}, v)]$ exists, then the situation is symmetric. Therefore, for balance of $[v, v^{-1}]$, the self-complemental biarc $[(v, v^{-1})]$ requires two biarcs entering $[v, v^{-1}]$ and the self-complemental biarc $[(v^{-1}, v)]$ requires two biarcs leaving $[v, v^{-1}]$. If both self-complemental arcs exist (e.g. both $[(ATA, TAT)]$ and $[(TAT, ATA)]$ for a binode $[ATA, TAT]$), then a biwalk can traverse them consecutively from e.g. $[v, v^{-1}]$ by traversing first $[(v, v^{-1})]$ and then $[(v^{-1}, v)]$, ending up in $[v, v^{-1}]$ again, such that the self-complemental arcs have a neutral contribution to the bi-imbalance. Resulting, the bi-imbalance of $[v, v^{-1}]$ is

$$bi_v = d^+(v) - d^-(v) + (\mathbb{1}_{(v, v^{-1}) \in E^+(v)} - \mathbb{1}_{(v^{-1}, v) \in E^-(v)}),$$

where $\mathbb{1}_P$ is 1 if the predicate P is true and 0 otherwise.

For self-complemental binodes $[v] \in V$, there is no concept of incoming or outgoing biarcs, since any biarc can be used to either enter or leave $[v]$ (see Fig. 1b for an example). Therefore, for balance, biarcs need to be paired arbitrarily, for which to be possible there needs to be an even amount of biarcs. If there is an odd amount, then there is one unpaired biarc, hence the bi-imbalance is 1. The following condition exerts this behaviour (using mod as the remainder operator):

$$bi_v = d^+(v) \bmod 2.$$

Finally, we include partial coverage to above bi-imbalance formulations by limiting the incoming and outgoing arcs to F . Further, to distinguish between self-complemental nodes and others, we denote the set of self-complemental nodes as $S \subseteq V$ and the set of binodes that are not self-complemental as $T := V \setminus S$. If self-complemental biarcs are included in the flow, then these alter the bi-imbalance by two, in the same way as they do in the equation of the bi-imbalance. We encode this on the left side of the equation. Then we get the following modified coverage constraint:

$$\begin{aligned} \forall v \in T : \sum_{e \in E^-(v)} x_e (1 + \mathbb{1}_{e=e^{-1}}) - \sum_{e \in E^+(v)} x_e (1 + \mathbb{1}_{e=e^{-1}}) \\ = d_F^+(v) - d_F^-(v) + (\mathbb{1}_{(v,v^{-1}) \in F^+(v)} - \mathbb{1}_{(v^{-1},v) \in F^-(v)}), \text{ and} \end{aligned} \quad (7)$$

$$\forall v \in S : \left(d_F^+(v) + \sum_{e \in E^+(v)} x_e \right) \bmod 2 = 0. \quad (8)$$

Valid bidirected flow

To adapt Edmonds' and Johnson's formulation to biwalks, we additionally need to ensure that the resulting flow yields a set of biarcs, i.e. that each arc has the same flow as its reverse complement:

$$\forall e \in E : x_e = x_{e^{-1}} \quad (9)$$

Adapted flow formulation

With the modifications above, we can adapt the formulation of Edmonds and Johnson to solve the bidirected partial-coverage Chinese postman problem. We define F to be the arcs in the original graph, and set $E := V \times V$. We further set $c_e = 1$ for $e \in F$ and $c_e = k - 1$ otherwise. Lastly we define S and T as above. Then we get the following modified formulation. Note that it is conceptually similar to that proposed in [56], however different because the basic definitions differ, and we further allow for special arcs of which only one needs to be covered.

$$\min \sum_{e \in E, e \text{ is canonical}} c_e x_e \text{ s.t.} \quad (5)$$

$$x_e \text{ are non-negative integers, and} \quad (3)$$

$$\sum_{e \in (E \setminus F)} x_e \geq 1, \text{ and} \quad (6)$$

$$\begin{aligned} \forall v \in T : \sum_{e \in E^-(v)} x_e (1 + \mathbb{1}_{e=e^{-1}}) - \sum_{e \in E^+(v)} x_e (1 + \mathbb{1}_{e=e^{-1}}) \\ = d_F^+(v) - d_F^-(v) + (\mathbb{1}_{(v, v^{-1}) \in F^+(v)} - \mathbb{1}_{(v^{-1}, v) \in F^-(v)}), \text{ and} \end{aligned} \quad (7)$$

$$\forall v \in S : \left(d_F^+(v) + \sum_{e \in E^+(v)} x_e \right) \bmod 2 = 0, \text{ and} \quad (8)$$

$$\forall e \in E : x_e = x_{e^{-1}} \quad (9)$$

In this min-cost integer flow formulation of the bidirected partial-coverage Chinese postman problem, analogue to the formulation of Edmonds and Johnson, sources and sinks are nodes with missing outgoing or incoming arcs, with demands matching the number of missing arcs in F . Our formulation would not be solvable for practical de-Bruijn graphs because inserting a quadratic amount of arcs into the graph is infeasible. However, most of the breaking arcs are not needed, since in a minimum solution they can only carry flow if they directly connect a source to a sink, by the following argument: Imagine a breaking arc that carries flow but is connected to a source or sink with at most one end. We can trace one unit of flow on the arc to a source and a sink, creating a path of flow one. By removing the flow from the path, and adding it to a breaking arc directly connecting the source to the sink, we get a valid flow. This flow has lower costs than the original, because it contains the same amount of breaking arcs, but a lower number of non-breaking arcs. This can be repeated until only breaking arcs that directly connect sources to sinks are left.

But even reducing the number of breaking arcs like this might not be enough if the graph contains too many sources and sinks. We therefore reduce the linear program to a min-cost matching instance, similar to Edmonds and Johnson.

Solving the min-cost integer flow formulation with min-cost matching

To solve the bidirected partial-coverage Chinese postman problem with min-cost matching, we observe that flow traverses the graph from a source to a sink only via min-cost paths, since all arcs have infinite capacity. Due to the existence of the breaking arcs with low costs ($k - 1$), we can further restrict the flow to use only paths of length at most $k - 2$ without affecting minimality. However, since we are also interested in a low number of strings in our minimum SPSS, we also allow paths of length $k - 1$. We can precompute these min-cost paths efficiently in parallel (see “[Efficient computation of many-to-many min-cost paths](#)” section below). Then it remains to decide which combination of min-cost paths and breaking arcs yield a minimum solution.

To simplify this problem, observe that the pairing of sources and sinks that are connected via breaking arcs does not matter. Any pairing uses the same amount of breaking arcs, and therefore has the same costs. It only matters that these nodes are not connected by a lower-cost path that does not use any breaking arcs, and that there is at least one breaking arc. As a result, we can ignore breaking arcs when searching a minimum

solution, and instead introduce costs for unmatched nodes. However, we still need to enforce that there is at least one pair of unmatched nodes. We do this using a special construction described below. Note though, that there can only be unmatched nodes if there are unbalanced binodes, i.e. the graph was not biEulerian in the beginning. However, if the graph is biEulerian already, the whole matching step can be left out, and instead a biEulerian circuit with one arbitrarily inserted breaking biarc can be returned. So we can safely assume here that the graph contains at least one pair of unbalanced binodes (it cannot contain a single unbalanced binode, see e.g. [47]).

We formulate a min-cost matching problem with penalty costs for unmatched nodes, which can be reduced to a min-cost perfect matching problem. For the construction of our undirected matching graph M we define the set of sources $A \subseteq T$ as all nodes with negative bi-imbalance, and the set of sinks $B \subseteq T$ as all nodes with positive bi-imbalance. Then we add $|bi_v|$ (absolute value of the bi-imbalance of v) copies of each node from A , B and S to M . Further, for each min-cost path from a node $a \in A \cup S$ to a node $b \in B \cup S$ we add an edge (undirected arc) from each copy of a to each copy of b in M with costs equal to the costs of the path. We ignore self loops at nodes in S since they do not alter the imbalance, and nodes in A and B cannot have self loops.

Then, to fulfil Condition (9) (valid bidirected flow) and to reduce the size of the matching problem, we merge all nodes and arcs with their reverse complement (the unmerged graph is built here to simplify our explanations, in our implementation we directly build the merged graph). This additionally results in self-complemental biarcs forming self-loops in the merged graph, thus making them not choosable by the matcher. But this is correct, as self-complemental biarcs alter the bi-imbalance of a binode by two, and therefore they can only be chosen by matching two different copies of the same binode in M .

Additionally, to fulfil Condition (6) (cover one bidirected arc), we add a pair of nodes u, w to M . We connect u and w to each node in M (but do not add an edge between u and w) and assign costs 0 to all those edges. This forces u and w to be matched to other nodes u', w' , which means that when biEulerising, the bi-imbances of u' and w' need to be fixed with at least one breaking arc.

Lastly, we assign each node other than u and w penalty costs of $(k - 1)/2$ for staying unmatched, as each pair of unmatched nodes produces costs $k - 1$ for using a breaking arc.

We reduce M to an instance of the min-cost perfect matching problem using the reduction described in [64]. For that we duplicate the graph, and add an edge with costs $k - 1$ between each node and its copy, but not between v and w and their respective copies. The costs for edges between a node and its copy are double the costs of keeping a node unmatched, because by choosing such an edge causes two nodes to stay unmatched.

After this reduction, we use Blossom V [65] to compute a solution. Since all nodes were doubled in the reduction, we actually get two solutions that might even differ, however both of them are minimum. We arbitrarily choose one of the solutions. This gives us a multiset of arcs that we complete with the breaking arcs required to balance the unmatched nodes to create a biEulerisation of the input graph. Following the approach

from Edmonds and Johnson, we find a biEulerian circuit in the resulting graph which is a solution to the bidirected partial-coverage Chinese postman problem as required.

Note that our matching formulation only optimises CL, but does not optimise SC. It indirectly optimises SC because it decreases CL by joining strings, which also decreases SC by one each time. However, when joining two strings would not alter CL, Blossom V may output both variants, with joining the strings and without, while staying optimal. It then chooses an arbitrary option.

Efficient computation of many-to-many min-cost paths

Apart from solving the matching problem, finding the min-cost paths between sources and sinks is the most computationally heavy part of our algorithm.

We solve it using Dijkstra's shortest path algorithm [66] in a one-to-many variant and execute it in parallel for all sources. To be efficient, we create a queue with blocks of source nodes, and the threads process one block at a time. A good block size balances between threads competing for access to the queue, and threads receiving a very imbalanced workload. Since our min-cost paths are short (at most $k - 1$ arcs), in most executions of Dijkstra's algorithm only a tiny fraction of the nodes in the graph are visited. But the standard variant of Dijkstra's algorithm wastefully allocates an array for all nodes to store their distance from the source node (the *distance array*). To save space, we instead use a hashmap, mapping from `node_index` to distance from source. This turned out to be faster than using a distance array, even if the distance array uses an epoch system to only do a full reset every 2^{32} queries. An epoch system stores a second value for each entry in the distance array indicating in what execution of Dijkstra's algorithm that value is valid. The execution counter gets incremented each execution, and only when it wraps around, the distance array is reset normally. As another optimisation, we abort the execution early when Dijkstra reaches costs greater than $k - 1$, since we are only interested in paths up to costs $k - 1$.

Finally, in our implementation, we do not compute the actual sequences of arcs of the paths. Instead of copying the path arcs when biEulerising the graph, we insert special dummy arcs with a length equal to the length of the path. When breaking the final biEulerian circuit, if there are no breaking arcs but dummy arcs, then we break at a longest dummy arc to produce a minimum solution. If there are neither breaking nor dummy arcs, we proceed as described above. Then, when reporting the final set of strings, we define `spell(\cdot)` to append the last ℓ characters of b when encountering a dummy biarc $[(a, b), (b^{-1}, a^{-1})]$ (or $[(a, b)]$) of length ℓ .

Efficient computation of the greedy heuristic

The greedy heuristic biEulerises the graph by greedily adding min-cost paths between unbalanced nodes, as opposed to choosing an optimal set via min-cost matching like our main algorithm. It then continues like the main algorithm, finding a biEulerian circuit, breaking it into walks and spelling out the strings.

To be efficient, the min-cost paths are again computed in parallel, and we apply all optimisations from “[Efficient computation of many-to-many min-cost paths](#)” section. The parallelism however poses a problem for the greedy computation: if a binode with one missing incoming biarc is reached by two min-cost paths in parallel, then if both

threads add their respective biarcs, we would overshoot the bi-imbalance of that binode. To prevent that, we introduce a lock for each node, and before inserting a biarc into the graph we lock all (up to) four incident nodes. By locking the nodes in order of their ids we ensure that no deadlock can occur. Since the number of threads is many orders of magnitude lower than the number of nodes, we assume that threads almost never need to wait for each other. In addition to the parallelisation, we abort Dijkstra's algorithm early when we have enough paths to fix the imbalance for the binode. This sometimes requires to execute Dijkstra's algorithm again if a potential sink node was used by a different thread in parallel. But again, since the number of threads is many orders of magnitude lower than the number of nodes, we assume that this case almost never occurs.

In practice, the greedy heuristic achieves better results in terms of SC than the optimal matchings algorithm (see the "Results" section). This is because the greedy heuristics always joins two strings if it does not alter CL, while the optimal algorithm does not, as explained in "Solving the min-cost integer flow formulation with min-cost matching" section.

Minimising string count

In the paper we studied SPSSes of minimum total length (minimum CL). In this section, we note that an SPSS with a minimum number of strings (minimum SC), and with no constraint on the total length, is also computable in polynomial time.

The high-level idea, ignoring reverse complements, is as follows. Given the arc-centric de Bruijn graph G , construct the directed acyclic graph G^* of strongly connected components (SCCs) of G . In G^* , every SCC is a node, and we have as many arcs between two SCCs as there are pairs of nodes in the two SCCs with an arc between them. Clearly, all arcs in a single SCC are coverable by a single walk. Moreover, for two SCCs connected by an arc, their two covering walks can be connected via this arc into a single walk covering all arcs of both SCCs. Thus, the minimum number of walks needed to cover all arcs of G^* (i.e. minimum SC SPSS) equals the minimum number of paths needed to cover all arcs of the acyclic graph G^* . This is a classic problem solvable in polynomial time with network flows (see e.g. [67] among many).

However, such an SPSS of minimum SC very likely has a large CL, because covering an SCC with a single walk might repeat quadratically many arcs, and connecting the covering walks of two adjacent SCCs might additionally require to repeat many arcs to reach the arc between them.

Experimental evaluation

We ran our experiments on a server running Linux with two 64-core AMD EPYC 7H12 processors with 2 logical cores per physical core, 1.96TiB RAM and an SSD. We downloaded the genomes of the model organisms from RefSeq [42]: *Caenorhabditis elegans* with accession GCF_000002985.6 [68], *Bombyx mori* with accession GCF_000151625.1 [69] and *Homo sapiens* with accession GCF_000001405.39 [70]. These are the same genomes as in [43], except that we downloaded HG38 from RefSeq for citability. The short reads were downloaded from the sequence read archive [71]: *Caenorhabditis elegans* with accession SRR14447868.1 [72], *Bombyx mori* with

accession DRR064025.1 [73] and *Homo sapiens* with accessions SRR2052337.1 to SRR2052425.1 [74].

We downloaded the 1102 *Neisseria gonorrhoeae* genomes from [75]. We downloaded the 616 *Streptococcus pneumoniae* genomes from the sequence read archive, using the accession codes provided in Table 1 in [76]. Up to here the pangenomes are retrieved in the same way as in [43]. We additionally used `grep` to select 3682 *Escherichia coli* genomes from GenBank [77] using the overview file ftp://ftp.ncbi.nlm.nih.gov/genomes/genbank/bacteria/assembly_summary.txt. The genomes are listed in Additional file 2. The ~309k salmonella genome sequences were downloaded from the EnteroBase database [78] in February 2022. The included filenames are in Additional file 3. The 2505x human pangenome is from the 1000 genomes project [79], created by downloading a variant of GRCh37 from ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/reference/phase2_reference_assembly_sequence/hs37d5.fa.gz and downloading variant files for chromosomes 1-22 from <http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/>. We then converted chromosomes 1-22 in the reference into 2505 sequences each using the tool `vcf2multialign` published in [80].

For querying the human read dataset, we used 3 million reads randomly drawn from the reads used to construct the dataset. For querying the *Salmonella* pangenome, we used 3 million randomly drawn short reads from 10 read data sets from the sequence read archive with the accessions listed in Additional file 4. For querying the human pangenome, we used 3 million randomly drawn short reads from one file of the sequence read archive with accession SRR2052337.1. This is one of the files from the human short read dataset described above. For querying the *E.coli* pangenome (in Additional file 1: Section 3) we used 30 short read data sets from the sequence read archive with the accessions listed in Additional file 5.

We used `snakemake` [81] and the `bioconda` software repository [82] to craft our experiment pipeline. The `tigs` were checked for correctness by checking the k -mer sets against unitigs. The `bifrost` queries in Additional file 1: Section 3 were checked for correctness by checking that the query results are equivalent for those with unitigs. The `SSHash` queries were not checked for correctness, as `SSHash` was modified by the author himself. Whenever we measured runtime of queries and builds for Additional file 1: Supplemental figure 5 (Performance with different amounts of threads), we only let a single experiment run, even if the experiment used only one core. When running the other builds we ran multiple processes at the same time, but never using more threads than the processor has physical cores (thus avoiding any effects introduced by sharing logical cores). When running a tool we copied its input to the SSD, and copied the results back to our main network storage, to prevent the network storage's varying workload to affect our runtime measurements. For experiments running on input reads or references as opposed to unitigs (`BCALM2`, `ProphAsm`), we copied the inputs to a RAID of HDDs instead, due to their large size. The copying was not part of the measurements. We made sure that the server is undisturbed, except that we monitored the experiment status and progress with `htop` and `less`. We limited each run to 256GiB of RAM per process, which prevented us from running `matchtigs` on larger inputs. Further, `ProphAsm` supports only $k \leq 32$, so it was not run for k larger than that.

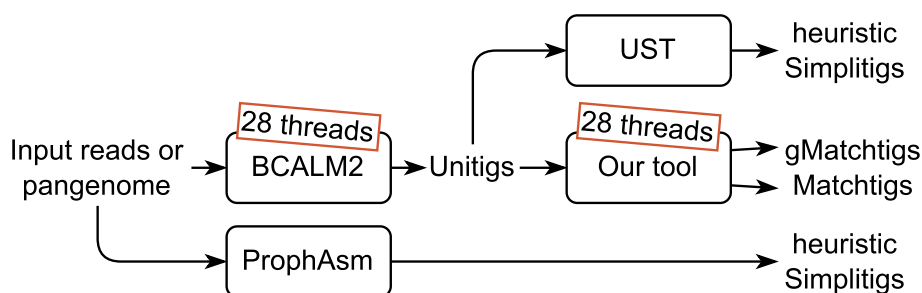


Fig. 5 Toolchain for computing various tigs. The DAG of tools run on the short read sets and pangenomes to compute different tigs

For an overview of our experiment pipeline for computing tigs, see Fig. 5. We run ProphAsm on the input data, as it was introduced to do [43]. All other tools require unitigs to be computed first. UST specifically requires unitigs computed by BCALM2, as BCALM2 adds additional annotations to the fasta unitig file. Our tool matchtigs also can make use of these annotations to speed up the construction of the arc-centric de Bruijn graph. On the human pangenome, BCALM2 crashed due to the input being too large. Hence we used Cuttlefish 2 [40] to compute unitigs, and since UST only runs on unitigs computed by BCALM2, we then ran ProphAsm to compute heuristic simplitigs.

For queries, we executed Bifrost or SShash-Lite on the different tigs. The Bifrost query command handles both building the index and executing the query, while SShash-Lite requires to run a separate command to build the index first.

We measure runtimes and memory using the command `/usr/bin/time -v`. The runtimes are reported as wall clock time, and the memory is reported as the maximum resident set size.

See Section *Availability of data and materials* for availability of our implementation and experiment code, which includes all the concrete commands we have used to execute our experiments.

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-023-02968-z>.

Additional file 1. Additional experiments and proof of optimality. Additional CL and SC data from our experiments, with varying k and min abundance. Also, the average unitig length and total unitig count is plotted. Performance measurements with varying amount of threads. A query experiment with Bifrost as well as an SShash-Lite experiment on a machine with focus on single-core performance. Proof of optimality for our algorithm.

Additional file 2. Accessions of the genomes used for the E. coli pangenome. A list of links to the genomes in fasta format on NCBI.

Additional file 3. File names of the genomes used for the Salmonella pangenome. A list of names of files downloaded from EnteroBase [78].

Additional file 4. Accessions of the reads used for the Salmonella query experiment. A list of sequence read archive accession numbers for the 10 sets of Salmonella short reads used for querying.

Additional file 5. Accessions of the reads used for the E. coli query experiment. A list of sequence read archive accession numbers for the 30 sets of E. coli short reads used for querying.

Additional file 6. Review history.

Acknowledgements

We are very grateful to Paul Medvedev and Amatur Rahman for helpful initial discussions on this problem. We further wish to thank the Finnish Computing Competence Infrastructure (FCCI) for supporting this project with computational and data storage resources. We also wish to thank Andrea Cracco for providing us with the ~ 309 kx Salmonella

pangenome. We further wish to thank the anonymous reviewers for their useful constructive feedback, which improved the presentation of the paper, the implementation and the experimental results. Finally we wish to thank the Rust community (<https://users.rust-lang.org>) for explanations about language-specific details of parallel implementations.

Peer review information

Andrew Cosgrove was the primary editor of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

Review history

The review history is available as Additional file 6.

Authors' contributions

SK and AIT formulated the problem. AIT and SK designed an optimal algorithm and SK implemented and evaluated its prototype. Thereafter, SS improved the algorithm's design to its current form and provided the final implementation as well as optimisations required to make it practically relevant. SS developed the greedy heuristic and implemented it. JA, SS and AIT designed the experiments and interpreted the results. SS performed all experiments (except for those on SSHash-Lite), and developed all further code published in the context of this work. SS wrote the manuscript under the supervision of the other authors, except for "k-mer-based short read queries" section. GEP created and described SSHash-Lite in "k-mer-based short read queries" section and ran the query experiments on the machine with focus on single-core performance. All authors reviewed and approved the final version of the manuscript.

Authors' twitter handles

Sebastian Schmidt @Sebasti66652811, Jarno Alanko @jonalanko, Giulio E. Pibiri @giulio_pibiri, and Alexandru I. Tomescu @AlexTomescu0.

Funding

Open Access funding provided by University of Helsinki including Helsinki University Central Hospital. This work was partially funded by the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No. 851093, SAFE BIO), and by the Academy of Finland (grants No. 322595, 328877). This work was also partially supported by the project MobiDataLab (EU H2020 RIA, grant agreement No. 101006879). Funding for this research has also been provided by the European Union's Horizon Europe research and innovation programme (EFRA project, Grant Agreement Number 101093026).

Availability of data and materials

matchtigs: <https://github.com/algbio/matchtigs>.

SSHash-Lite: <https://github.com/jermp/sshash-lite>.

The implementation of the matchtigs and greedy matchtigs algorithms is available on GitHub [83]. The name of the project is *matchtigs*. It is platform independent, and can be compiled locally or installed from bioconda as described in the README of the project. It is licensed under the 2-clause BSD license. The version used for our experiments is available at [84], and the implementation together with all code to reproduce the experiments is available at [85]. The experiment code is also licensed under the 2-clause BSD license. SSHash-Lite is available on GitHub [86] and licensed under the MIT license. The version used for our experiments is available at [87] and licensed under the MIT license.

The genomes of the model organisms are: *Caenorhabditis elegans* [68], *Bombyx mori* [69] and *Homo sapiens* [70]. The short reads of the model organisms are: *Caenorhabditis elegans* [72], *Bombyx mori* [73] and *Homo sapiens* [74].

The 616x *Streptococcus pneumoniae* pangenome is available in the sequence read archive [71], using the accession codes provided in Table 1 in [76]. The 1102x *Neisseria gonorrhoeae* pangenome is from [75]. The 3682x *Escherichia coli* pangenome is available from GenBank [77] using the accessions in Additional file 2. The ~309kx *Salmonella* pangenome is created by downloading all *Salmonella* genomes from the Enterobase database [78] in February 2022. The included filenames are in Additional file 3. The 2505x human pangenome is from the 1000 genomes project [79].

The query for the *Salmonella* pangenome are 3 million randomly drawn short reads from 10 *Salmonella* read data sets from the sequence read archive with the accessions listed in Additional file 4. The query for the human pangenome are 3 million randomly drawn short reads from [74] restricted to file SRR2052337.1. The query for the *E. coli* pangenome (in Additional file 1: Section 3) are 30 short read data sets from the sequence read archive with the accessions listed in Additional file 5.

Declarations

Ethics approval and consent to participate

This study only uses publicly available datasets, hence an ethics approval or consent to participate is not required.

Competing interests

The authors declare that they have no competing interests.

Received: 16 December 2021 Accepted: 10 May 2023

Published online: 09 June 2023

References

- Zielezinski A, Vinga S, Almeida J, Karlowski WM. Alignment-free sequence comparison: benefits, applications, and tools. *Genome Biol.* 2017;18(1):1–17.
- Zielezinski A, Girgis HZ, Bernard G, Leimeister C-A, Tang K, Dencker T, Lau AK, Röhling S, Choi JJ, Waterman MS, et al. Benchmarking of alignment-free sequence comparison methods. *Genome Biol.* 2019;20(1):1–18.
- Luhmann N, Holley G, Achtman M. Blastfrost: fast querying of 100,000 s of bacterial genomes in bifrost graphs. *Genome Biol.* 2021;22(1):1–15.
- Iqbal Z, Caccamo M, Turner I, Flicek P, McVean G. De novo assembly and genotyping of variants using colored de Bruijn graphs. *Nat Genet.* 2012;44(2):226–32.
- Nordström KJ, Albani MC, James GV, Gutjahr C, Hartwig B, Turck F, Paszkowski U, Coupland G, Schneeberger K. Mutation identification by direct comparison of whole-genome sequencing data from mutant and wild-type individuals using k-mers. *Nat Biotechnol.* 2013;31(4):325–30.
- Bradley P, Gordon NC, Walker TM, Dunn L, Heys S, Huang B, Earle S, Pankhurst LJ, Anson L, De Cesare M, et al. Rapid antibiotic-resistance predictions from genome sequence data for staphylococcus aureus and mycobacterium tuberculosis. *Nat Commun.* 2015;6(1):1–15.
- Shajii A, Yorukoglu D, William YuY, Berger B. Fast genotyping of known snps through approximate k-mer matching. *Bioinformatics.* 2016;32(17):538–44.
- Sun C, Medvedev P. Toward fast and accurate snp genotyping from whole genome sequencing data for bedside diagnostics. *Bioinformatics.* 2019;35(3):415–20.
- Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic rna-seq quantification. *Nat Biotechnol.* 2016;34(5):525–7.
- Ames SK, Hysom DA, Gardner SN, Lloyd GS, Gokhale MB, Allen JE. Scalable metagenomic taxonomy classification using a reference genome database. *Bioinformatics.* 2013;29(18):2253–60.
- Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biol.* 2014;15(3):1–12.
- Břinda K, Salikhov K, Pignotti S, Kucherov G. Prophyle: a phylogeny-based metagenomic classifier using the burrows-wheeler transform. Poster at HiTSeq 2017. 2017.
- Corvelo A, Clarke WE, Robine N, Zody MC. taxmaps: comprehensive and highly accurate taxonomic classification of short-read data in reasonable time. *Genome Res.* 2018;28(5):751–8.
- Simon HY, Siddle KJ, Park DJ, Sabeti PC. Benchmarking metagenomics tools for taxonomic classification. *Cell.* 2019;178(4):779–94.
- Sirén J. Indexing variation graphs. In: 2017 Proceedings of the Nineteenth Workshop on Algorithm Engineering and Experiments (ALENEX). SIAM; 2017. pp. 13–27.
- Garrison E, Sirén J, Novak AM, Hickey G, Eizenga JM, Dawson ET, Jones W, Garg S, Markello C, Lin MF, et al. Variation graph toolkit improves read mapping by representing genetic variation in the reference. *Nat Biotechnol.* 2018;36(9):875–9.
- Benoit G. Simka: fast kmer-based method for estimating the similarity between numerous metagenomic datasets. In: RCAM. Le Chesnay Cedex: Inria Domaine de Voluceau Rocquencourt; 2015.
- David S, Mentasti M, Tewolde R, Aslett M, Harris SR, Afshar B, Underwood A, Fry NK, Parkhill J, Harrison TG. Evaluation of an optimal epidemiological typing scheme for legionella pneumophila with whole-genome sequence data using validation guidelines. *J Clin Microbiol.* 2016;54(8):2135–48.
- Chattaway MA, Schaefer U, Tewolde R, Dallman TJ, Jenkins C. Identification of escherichia coli and shigella species from whole-genome sequences. *J Clin Microbiol.* 2017;55(2):616–23.
- Klausen PT, Aarestrup FM, Lund O. Rapid and precise alignment of raw reads against redundant databases with KMA. *BMC Bioinformatics.* 2018;19(1):1–8.
- Pandey P, Almodaresi F, Bender MA, Ferdman M, Johnson R, Patro R. Mantis: A fast, small, and exact large-scale sequence-search index. *Cell Syst.* 2018;7(2):201–7.
- Marchet C, Kerbirou M, Limasset A. Indexing De Bruijn graphs with minimizers. In: Recomb-Seq 2019-9th RECOMB Satellite Workshop on Massively Parallel Sequencing. Le Chesnay Cedex: Inria Domaine de Voluceau Rocquencourt; 2019. pp. 1–16.
- Holley G, Melsted P. Bifrost: highly parallel construction and indexing of colored and compacted de Bruijn graphs. *Genome Biol.* 2020;21(1):1–20.
- Pevzner PA. I-Tuple DNA sequencing: computer analysis. *J Biomol Struct Dyn.* 1989;7(1):63–73.
- Idury RM, Waterman MS. A new algorithm for DNA sequence assembly. *J Comput Biol.* 1995;2(2):291–306.
- Pevzner PA, Tang H, Waterman MS. An Eulerian path approach to DNA fragment assembly. *Proc Natl Acad Sci.* 2001;98(17):9748–53. <https://doi.org/10.1073/pnas.171285098>.
- Chaisson MJ, Pevzner PA. Short read fragment assembly of bacterial genomes. *Genome Res.* 2008;18(2):324–30.
- Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJ, Birol I. Abyss: a parallel assembler for short read sequence data. *Genome Res.* 2009;19(6):1117–23.
- Li R, Zhu H, Ruan J, Qian W, Fang X, Shi Z, Li Y, Li S, Shan G, Kristiansen K, et al. De novo assembly of human genomes with massively parallel short read sequencing. *Genome Res.* 2010;20(2):265–72.
- Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko SI, Pham S, Pribelski AD, Pyskhin AV, Sirotkin AV, Vyahhi N, Tesler G, Alekseyev MA, Pevzner PA. SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing. *J Comput Biol.* 2012;19(5):455. <https://doi.org/10.1089/cmb.2012.0021>.
- Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Pan Q, Liu Y, et al. Soapdenovo2: an empirically improved memory-efficient short-read de novo assembler. *Gigascience.* 2012;1(1):2047–217.
- Chikhi R, Rizk G. Space-efficient and exact de Bruijn graph representation based on a Bloom filter. *Algorithm Mol Biol.* 2013;8(1):22.
- Chikhi R, Limasset A, Medvedev P. Compacting de Bruijn graphs from sequencing data quickly and in low memory. *Bioinformatics.* 2016;32(12):201–8.

34. Jackman SD, Vandervalk BP, Mohamadi H, Chu J, Yeo S, Hammond SA, Jahesh G, Khan H, Coombe L, Warren RL, Birol I. ABySS 2.0: resource-efficient assembly of large genomes using a Bloom filter. *Genome Res.* 2017;27(5):768–777. <https://doi.org/10.1101/gr.214346.116>.
35. Ruan J, Li H. Fast and accurate long-read assembly with wtdbg2. *Nat Methods.* 2020;17(2):155–8.
36. Tomescu AI, Medvedev P. Safe and Complete Contig Assembly Through Omnitigs. *J Comput Biol.* 2017;24(6):590–602. <https://doi.org/10.1089/cmb.2016.0141>.
37. Acosta NO, Mäkinen V, Tomescu AI. A safe and complete algorithm for metagenomic assembly. *Algorithm Mol Biol.* 2018;13(1):1–12.
38. Cairo M, Khan S, Rizzi R, Schmidt S, Tomescu AI, Zironelli EC. The hydrostructure: a universal framework for safe and complete algorithms for genome assembly. 2020. arXiv preprint [arXiv:2011.12635](https://arxiv.org/abs/2011.12635).
39. Kececioglu JD, Myers EW. Combinatorial algorithms for DNA sequence assembly. *Algorithmica.* 1995;13(1):7–51.
40. Khan J, Kokot M, Deorowicz S, Patro R. Scalable, ultra-fast, and low-memory construction of compacted de bruijn graphs with Cuttlefish 2. *Genome Biol.* 2022;23(1):1–32.
41. Cracco A, Tomescu AI. Extremely-fast construction and querying of compacted and colored de Bruijn graphs with GGCA. *bioRxiv.* 2022. <https://doi.org/10.1101/2022.10.24.513174>. <https://www.biorxiv.org/content/early/2022/10/25/2022.10.24.513174.full.pdf>
42. O’Leary NA, Wright MW, Brister JR, Ciuffo S, Haddad D, McVeigh R, Rajput B, Robbertse B, Smith-White B, Ako-Adjei D, et al. Reference sequence (refseq) database at ncbi: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res.* 2016;44(D1):733–45.
43. Břinda K, Baym M, Kucherov G. Simplitigs as an efficient and scalable representation of de Bruijn graphs. *Genome Biol.* 2021;22(1):1–24.
44. Rahman A, Medvedev P. Representation of k-mer sets using spectrum-preserving string sets. *J Comput Biol.* 2021;28(4):381–94.
45. Pibiri GE. Sparse and skew hashing of k-mers. *Bioinformatics.* 2022;38(Supplement_1):185–194. <https://doi.org/10.1093/bioinformatics/btac245>.
46. Li H. Aligning sequence reads, clone sequences and assembly contigs with bwa-mem. 2013. arXiv preprint [arXiv:1303.3997](https://arxiv.org/abs/1303.3997).
47. Schmidt S, Alanko JN. Eulertigs: Minimum Plain Text Representation of k-mer Sets Without Repetitions in Linear Time. In: Boucher C, Rahmann S, editors. 22nd International Workshop on Algorithms in Bioinformatics (WABI 2022). Leibniz International Proceedings in Informatics (LIPIcs), vol 242. pp. 1–21. Dagstuhl: Schloss Dagstuhl – Leibniz-Zentrum für Informatik; 2022. <https://doi.org/10.4230/LIPIcs.WABI.2022.2>. <https://drops.dagstuhl.de/opus/volltexte/2022/17036>.
48. Dufresne Y, Lemane T, Marijon P, Peterlongo P, Rahman A, Kokot M, Medvedev P, Deorowicz S, Chikhi R. The k-mer file format: a standardized and compact disk representation of sets of k-mers. *Bioinformatics.* 2022;38(18):4423–5.
49. Fan J, Khan J, Pibiri GE, Patro R. Spectrum preserving tilings enable sparse and modular reference indexing. *bioRxiv.* 2022. <https://doi.org/10.1101/2022.10.27.513881>. <https://www.biorxiv.org/content/early/2022/10/28/2022.10.27.513881.full.pdf>.
50. Kitaya K, Shibuya T. Compression of Multiple k-Mer Sets by Iterative SPSS Decomposition. In: Carbone A, El-Kebir M, editors. 21st International Workshop on Algorithms in Bioinformatics (WABI 2021). Leibniz International Proceedings in Informatics (LIPIcs), vol 201. Dagstuhl: Schloss Dagstuhl – Leibniz-Zentrum für Informatik; 2021. pp. 12–11217. <https://doi.org/10.4230/LIPIcs.WABI.2021.12>. <https://drops.dagstuhl.de/opus/volltexte/2021/14365>.
51. Marchet C, Iqbal Z, Gautheret D, Salson M, Chikhi R. REINDEER: efficient indexing of k-mer presence and abundance in sequencing datasets. *Bioinformatics.* 36(Supplement_1):177–185. 2020. <https://doi.org/10.1093/bioinformatics/btaa487>. https://academic.oup.com/bioinformatics/article-pdf/36/Supplement_1/177/33860751/btaa487.pdf.
52. Rahman A, Chikhi R, Medvedev P. Disk compression of k-mer sets. *Algorithm Mol Biol.* 2021;16(1):1–14.
53. Kwan M-k. Graphic programming using odd or even points. *Chin Math.* 1962;1:273–7.
54. Edmonds J, Johnson EL. Matching, euler tours and the chinese postman. *Math Program.* 1973;5(1):88–124.
55. Kundeti V, Rajasekaran S, Dinh H. An efficient algorithm for chinese postman walk on bi-directed de bruijn graphs. In: Wu W, Daescu O, editors. *Combinatorial Optimization and Applications*. Berlin, Heidelberg: Springer; 2010. p. 184–96.
56. Medvedev P, Georgiou K, Myers G, Brudno M. Computability of models for sequence assembly. In: Giancarlo R, Hannenhalli S, editors. *Algorithms in Bioinformatics, 7th International Workshop, WABI 2007, Philadelphia, PA, USA, September 8-9, 2007, Proceedings*. Lecture Notes in Computer Science, vol 4645. Berlin, Heidelberg: Springer; 2007. pp. 289–301. https://doi.org/10.1007/978-3-540-74126-8_27.
57. Pibiri GE. On Weighted k-mer Dictionaries. In: Boucher C, Rahmann S, editors. 22nd International Workshop on Algorithms in Bioinformatics (WABI 2022). Leibniz International Proceedings in Informatics (LIPIcs), vol 242. Dagstuhl: Schloss Dagstuhl – Leibniz-Zentrum für Informatik; 2022. pp. 1–20. <https://doi.org/10.4230/LIPIcs.WABI.2022.9>. <https://drops.dagstuhl.de/opus/volltexte/2022/17043>
58. Pibiri GE, Trani R. PTHash: Revisiting FCH Minimal Perfect Hashing. In: *The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. New York: Association for Computing Machinery; 2021. pp. 1339–1348.
59. Roberts M, Hayes W, Hunt BR, Mount SM, Yorke JA. Reducing storage requirements for biological sequence comparison. *Bioinformatics.* 2004;20(18):3363–9.
60. Lenstra JK, Kan AR. Complexity of vehicle routing and scheduling problems. *Networks.* 1981;11(2):221–7.
61. Edmonds J, Karp RM. Theoretical improvements in algorithmic efficiency for network flow problems. *J ACM (JACM).* 1972;19(2):248–64.
62. Christofides N, Campos V, Corberán A, Mota E. In: Gallo G, Sandi C, editors. *An algorithm for the Rural Postman problem on a directed graph*. Berlin, Heidelberg: Springer; 1986. pp. 155–166. <https://doi.org/10.1007/BFb0121091>.
63. Even S. *Graph Algorithms*. Rockville: Computer Science Press; 1979.
64. Schäfer G. *Weighted matchings in general graphs*. Master’s thesis, Saarland University; 2000.

65. Kolmogorov V, Blossom V: a new implementation of a minimum cost perfect matching algorithm. *Math Program Comput.* 2009;1(1):43–67.
66. Dijkstra EW. A note on two problems in connexion with graphs. *Numer Math.* 1959;1(1):269–71.
67. Cáceres M, Cairo M, Mumey B, Rizzi R, Tomescu AI. Sparsifying, shrinking and splicing for minimum path cover in parameterized linear time. 2021. arXiv preprint [arXiv:2107.05717](https://arxiv.org/abs/2107.05717). To appear in the Proceedings of the 2022 ACM-SIAM Symposium on Discrete Algorithms (SODA 2022).
68. C. elegans Sequencing Consortium. *Caenorhabditis elegans* Bristol N2. 2013. https://www.ncbi.nlm.nih.gov/assembly/GCF_000002985.6/. Accessed 18 Apr 2023.
69. The international silkworm genome sequencing consortium. *Bombyx mori* p50T (= Dazao). 2008. https://www.ncbi.nlm.nih.gov/assembly/GCF_000151625.1/. Accessed 18 Apr 2023.
70. Genome Reference Consortium. Genome Reference Consortium Human Build 38 patch release 13 (GRCh38.p13). 2019. https://www.ncbi.nlm.nih.gov/assembly/GCF_000001405.39/. Accessed 18 Apr 2023.
71. Leinonen R, Sugawara H, Shumway M. The sequence read archive. *Nucleic Acids Res.* 2010;39(suppl_1):19–21.
72. Institute of Genetics and Developmental Biology. Deep sequencing of *Caenorhabditis elegans* with transgenerational UPPrmt. 2021. <https://www.ncbi.nlm.nih.gov/sra/?term=SRR14447868>. Accessed 18 Apr 2023.
73. University of Tokyo - Graduate School of Agricultural and Life Sciences (UT-GALS). Illumina HiSeq 2500 paired end sequencing of SAMD00054089. 2016. <https://www.ncbi.nlm.nih.gov/sra/?term=DRR064025>. Accessed 18 Apr 2023.
74. NCBI. NIST Genome in a Bottle, 300X sequencing of HG001 (NA12878)-131219_D00360_005_BH814YADXX. 2015. <https://www.ncbi.nlm.nih.gov/sra/?term=SRR2052337> to <https://www.ncbi.nlm.nih.gov/sra/?term=SRR2052425>. Accessed 18 Apr 2023.
75. Grad Y. Data for “Genomic Epidemiology of Gonococcal Resistance to Extended-Spectrum Cephalosporins, Macrolides, and Fluoroquinolones in the United States, 2000–2013”. Zenodo. 2019. <https://doi.org/10.5281/zenodo.2618836>.
76. Croucher NJ, Finkelstein JA, Pelton SI, Parkhill J, Bentley SD, Lipsitch M, Hanage WP. Population genomic datasets describing the post-vaccine evolutionary epidemiology of *Streptococcus pneumoniae*. *Sci Data.* 2015;2(1):1–9.
77. Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. GenBank. *Nucleic Acids Research.* 2012;41(D1):36–42. <https://doi.org/10.1093/nar/gks1195>. <https://academic.oup.com/nar/article-pdf/41/D1/D36/3680750/gks1195.pdf>
78. Zhou Z, Alikhan N-F, Mohamed K, Fan Y, Achtman M. The user’s guide to comparative genomics with Enterobase, including case studies on transmissions of micro-clades of *Salmonella*, the phylogeny of ancient and modern *Yersinia pestis* genomes, and the core genomic diversity of all *Escherichia*. *bioRxiv.* 2019. <https://doi.org/10.1101/613554>. <https://www.biorxiv.org/content/early/2019/11/25/613554.full.pdf>.
79. Consortium GP, et al. A global reference for human genetic variation. *Nature.* 2015;526(7571):68.
80. Norri T, Cazaux B, Dönges S, Valenzuela D, Mäkinen V. Founder reconstruction enables scalable and seamless pangenomic analysis. *Bioinformatics.* 2021;37(24):4611–9.
81. Köster J, Rahmann S. Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics.* 2012;28(19):2520–2.
82. Grüning B, Dale R, Sjödin A, Chapman BA, Rowe J, Tomkins-Tinch CH, Valieris R, Köster J. Bioconda: sustainable and comprehensive software distribution for the life sciences. *Nat Methods.* 2018;15(7):475–6.
83. Schmidt S. Matchtigs. GitHub. 2022. <https://github.com/algbio/matchtigs>. Accessed 18 Apr 2023.
84. Matchtigs Schmidt S. Zenodo. 2022. <https://doi.org/10.5281/zenodo.7371184>.
85. Schmidt S. Matchtigs experiments. Zenodo. 2022. <https://doi.org/10.5281/zenodo.7275990>.
86. Pibiri GE. SSHash-Lite. GitHub. 2022. <https://github.com/jermp/sshash-lite>. Accessed 18 Apr 2023.
87. Pibiri GE. SSHash-Lite Zenodo. 2022. <https://doi.org/10.5281/zenodo.7277145>.

Publisher’s Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

