

1 Compressed Weighted de Bruijn Graphs

2 **Giuseppe F. Italiano** 

3 Luiss University, Rome, Italy

4 Erable, INRIA Grenoble Rhône-Alpes, France

5 gitaliano@luiss.it

6 **Nicola Prezza** 

7 DAIS, Ca' Foscari University of Venice, Italy

8 nicola.prezza@unive.it

9 **Blerina Sinaimer** 

10 Luiss University, Rome, Italy

11 Erable, INRIA Grenoble Rhône-Alpes, France

12 bsinaimer@luiss.it

13 **Rossano Venturini** 

14 Dipartimento di Informatica, Università di Pisa, Pisa, Italy

15 rossano.venturini@unipi.it

16 — Abstract —

17 We propose a new compressed representation for weighted de Bruijn graphs, which is based on the
18 idea of delta-encoding the variations of k -mer abundances on a spanning branching of the graph.
19 Our new data structure is likely to be of practical value: to give an idea, when combined with the
20 compressed *BOSS* de Bruijn graph representation, it encodes the weighted de Bruijn graph of a
21 16x-covered DNA read-set (60M distinct k -mers, $k = 28$) within 4.15 bits per distinct k -mer and can
22 answer abundance queries in about 60 microseconds on a standard machine. In contrast, state of
23 the art tools declare a space usage of at least 30 bits per distinct k -mer for the same task, which
24 is confirmed by our experiments. As a by-product of our new data structure, we exhibit efficient
25 compressed data structures for answering partial sums on edge-weighted trees, which might be of
26 independent interest.

27 **Availability and implementation:** The code is written in C++ and is available at <https://github.com/nicolaprezza/cw-dBg>.

29 **2012 ACM Subject Classification** Theory of computation → Data compression; Theory of compu-
30 tation → Data structures design and analysis

31 **Keywords and phrases** weighted de Bruijn graphs, k -mer annotation, compressed data structures,
32 partial sums

33 **Digital Object Identifier** 10.4230/LIPIcs.CVIT.2016.23

34 **Funding** *Giuseppe F. Italiano*: Partially supported by MUR, the Italian Ministry for University and
35 Research, under PRIN Project AHeAD (Efficient Algorithms for HARnessing Networked Data).

36 *Rossano Venturini*: Partially supported by MUR, the Italian Ministry for University and Research,
37 under PRIN 2017 Project 2017K7XPAN *Algorithms, Data Structures and Combinatorics for Machine*
38 *Learning* and Università di Pisa Project PRA_2020-2021_26.

39 **1 Introduction**

40 A DNA resequencing experiment consists of reconstructing the nucleotide sequence of large
41 collections of short DNA fragments sampled from a reference genome [16, 18]. Depending
42 on the genome's length and on the number of sampled fragments, each genome position is
43 typically covered several times (up to a few hundred on average) by different fragments. One
44 of the most common (lossy) representations of such a dataset is a de Bruijn graph of order
45 k (k is usually chosen around 30) introduced in [31] and used by the majority of genome



© Giuseppe F. Italiano, Nicola Prezza, Blerina Sinaimer, Rossano Venturini;
licensed under Creative Commons License CC-BY

42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:16

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

46 and transcriptome assemblers [2, 24, 20, 37]: this is a combinatorial structure storing all
 47 k -mers occurring in the DNA fragments, and connecting two k -mers with an edge when their
 48 length- $(k + 1)$ concatenation also occurs in the dataset. While such a representation is already
 49 useful to discover important features of the sequenced genome (for example, single-point
 50 mutations correspond to bubbles in De Bruijn graphs, see e.g., [35, 36]), a much more useful
 51 representation should also record the *abundance* of each k -mer, i.e., how often it appears in
 52 the dataset. This augmented structure takes the name *weighted de Bruijn graph* and has
 53 many applications in both genome and transcriptome analysis. For instance, in *de novo*
 54 genome projects it can be used for error correction [22, 37] or to infer genome characteristics
 55 as repeat structures or rate of heterozygosity [21]. In transcriptome projects the role of
 56 k -mer abundances is even more important as the abundance of a gene does not only reflect
 57 its copy-number in the genome, but also and mostly its expression level. Hence, k -mer
 58 abundances are crucial for quantifying the expression level of transcripts [26, 30] but also for
 59 distinguishing SNPs from sequencing errors, or between different types of alternative splicing
 60 events [17, 23].

61 Due to the massive size of sequencing data, the main challenge is to design data structures
 62 for weighted de Bruijn graphs, which can scale up to very large instances such as genomes,
 63 genome collections or metagenomics datasets (see for example [24]). In this paper, we tackle
 64 exactly this problem and present new space-efficient data structures for storing weighted de
 65 Bruijn graphs.

66 1.1 State of the art

67 As far as (unweighted) de Bruijn graphs are concerned, Bowe et al. in [3] presented a
 68 space-efficient data structure — *BOSS* in the following — based on the concept of *prefix*
 69 *sorting*. Their idea is to sort co-lexicographically the k -mers appearing in the dataset and to
 70 append their outgoing labels to a sequence. It turns out that this sequence is a generalization
 71 to de Bruijn graphs of the celebrated Burrows-Wheeler transform [4], the text permutation
 72 at the core of the FM-index [12]. In its most basic form, the BOSS representation requires
 73 just 4 bits per edge to be stored and supports fast navigation queries on the de Bruijn graph.

74 It would be highly desirable to achieve the same rate of space efficiency also for weighted
 75 de Bruijn graphs. A naive strategy to extend the BOSS representation to weighted de Bruijn
 76 graphs could be to associate to each node a counter explicitly storing k -mer abundances.
 77 The large variance of these counters (that is, difference between the largest and mean
 78 abundance), however, requires some compression strategy in order to save space with respect
 79 to a straightforward uncompressed solution.

80 Although many data structures already present in the literature could support k -mer
 81 quantification, their representations are not space efficient [24]. For instance an extension
 82 of Bloom filters, known as counting Bloom filters [19], allocates a fixed number of bits per
 83 k -mer to store its count. This is clearly not space efficient for datasets where most of the
 84 k -mers have low abundance. To deal with this, variable-length counters have been proposed.
 85 Among them the counting quotient filter has been introduced in [28] and based on it two
 86 tools have been proposed: deBGR [27] and Squeakr [29]. Both tools are *approximate* in the
 87 sense that the returned counters could be wrong with low probability, but can also store exact
 88 counts at the price of a higher space usage. On a RNA-seq Human experiment composed of
 89 1.4 billion distinct k -mers ($k = 28$) with average abundance per k -mer equal to 27, deBGR
 90 [27] reports a space usage of 37 bit/ k -mer, while Squeakr [29] uses 79.8 bit/ k -mer. Another
 91 experiment on k -mer counting tools [29] reports that, on a dataset composed of 6.6 billion
 92 distinct k -mers ($k = 28$) with average abundance per k -mer equal to 14.9, the RAM usage of

93 Squeakr [29] is of 77.8 bit/ k -mer, the RAM usage of KMC2 [8] (signature-based) is of 139
 94 bit/ k -mer, and the RAM usage of Jellyfish2 [25] (hash-based) is of 80 bit/ k -mer.

95 Finally, in [33] an approximate data structure called Set-Min sketch has been proposed.
 96 This method takes advantage of the power-law distribution of the k -mer counts to reduce
 97 the error rate of the returned results. Set-Min sketch is implemented in the tool `fress` which
 98 achieves better memory consumption but does not guarantee exact counts. For a dataset
 99 of 6.6 billion distinct k -mers ($k = 28$) with average abundance per k -mer equal to 14.9, the
 100 declared RAM usage of `fress` is 16.9 bit/ k -mer with an error rate of 0.01.

101 1.2 Our contribution.

102 In this paper we propose a new representation for weighted de Bruijn graphs, which is based
 103 on the idea of delta-encoding the abundance variations on a spanning branching of the graph.
 104 To show the usefulness of our representation we test it on different real datasets obtaining
 105 each time a significant space improvement with respect to the state of the art. As an example,
 106 when combined with the compressed BOSS de Bruijn graph representation [3], our new data
 107 structure stores the weighted de Bruijn graph of a 16x-covered DNA read-set (60M distinct
 108 k -mers, $k = 28$) within 4.15 bits per distinct k -mer: just 1.4 bits per k -mer on top of the
 109 BOSS representation¹ [3].

110 Our data structure does not make any assumption on the distribution of the k -mer counts,
 111 which makes it appropriate to any dataset for both DNA and RNA. As a by-product of our
 112 data structure, we exhibit efficient compressed data structures for answering partial sums on
 113 edge-weighted trees, which might be of independent interest.

114 2 Notation

115 Logarithms are to the base 2. To simplify notation, we take $\log 1 = 1$. Notation $[n]$
 116 denotes the set $\{1, 2, \dots, n\}$. Throughout the paper, we work with the DNA alphabet
 117 $\Sigma = \{A, C, G, T\}$. A k -mer is a string from Σ^k . If $w \in \Sigma^*$ is a string, then $w[i, j]$ denotes
 118 substring $w[i] \cdot w[i+1] \cdots w[j]$, where the symbol "." represents the concatenation of characters.

119 **de Bruijn graphs** A de Bruijn graph is a directed graph $G = (\mathcal{V}, \mathcal{E})$ whose nodes are in
 120 bijection with the k -mers appearing in the dataset. To distinguish between k -mers and the
 121 corresponding nodes of G , we use the notation \hat{s} to denote the node of G corresponding to k -
 122 mer $s \in \Sigma^k$. Let $s, s' \in \Sigma^k$ be two k -mers. We say that s is adjacent to s' if $s[2, k] = s'[1, k-1]$
 123 and the $(k+1)$ -mer $s \cdot s'[k]$ occurs in the dataset. The directed edges of G are in bijection with
 124 all pairs of adjacent k -mers: (\hat{s}, \hat{s}') belongs to \mathcal{E} if and only if s is adjacent to s' . A *weighted*
 125 de Bruijn graph associates an integer positive weight $c(\hat{s})$ to each of its nodes. Weight $c(\hat{s})$
 126 (equivalently, $c(s)$ when we refer to the corresponding k -mer) is called the *abundance* of s
 127 and corresponds to the number of times that k -mer s appears in the dataset.

128 3 Compressed Weighted de Bruijn Graphs

129 To compress the weights (abundances) of a de Bruijn graph, one could exploit the following
 130 observation: since consecutive genomic positions generate adjacent k -mers, weights of adjacent
 131 nodes in the de Bruijn graph are likely to be very similar. Let $s, s' \in \Sigma^k$ be such that s is

¹ Note that we divide the total space by the number of *distinct* k -mers rather than the number of de Bruijn graph edges as done in [3]. More details on this in Section 4.

132 adjacent to s' . Recall that with $c(s)$ we denote the k -mer's abundance and we extend this
 133 notation to the nodes of G as $c(\hat{s}) = c(s)$. Let (\hat{s}, \hat{s}') be the corresponding edge of G , and let
 134 $w(\hat{s}, \hat{s}') = c(\hat{s}) - c(\hat{s}')$ be the difference between the weights of \hat{s} and \hat{s}' . Note that $w(\hat{s}, \hat{s}')$
 135 might be negative.

136 A first idea could be to store (i) the compressed integer $w(\hat{s}, \hat{s}')$ on each edge, and (ii)
 137 the explicit value $c(\hat{u})$ whenever \hat{u} has in-degree equal to 0. In this case, we say that \hat{u} is
 138 a *sampled node*. At this point, one could retrieve $c(\hat{s})$ also for any non-sampled node \hat{s} by
 139 summing $c(\hat{u})$ to the weights of the edges on a path from \hat{s} to \hat{u} , where \hat{u} is a sampled node.
 140 This solution has two main drawbacks: it is not time-efficient (in the worst case we need to
 141 visit the whole graph G to compute one single weight) and it stores one integer per edge,
 142 whereas the original graph contained one integer per node (the abundance).

143 In fact, all we need is a *spanning branching*² $T = (V, E)$ of G . The idea is to store (i) the
 144 compressed weight $w(\hat{s}, \hat{s}')$ only for the edges of T and (ii) the weights $c(\hat{r}_1), \dots, c(\hat{r}_t)$ for
 145 the t roots $\hat{r}_1, \dots, \hat{r}_t$ of T . To simplify the description, in the following we will assume that
 146 $c(\hat{r}_i) = 0$ for $1 \leq i \leq t$. This information is sufficient to reconstruct the weight of each node
 147 of T , but still leaves us with a few issues:

- 148 (1) We would like to design a fast data structure to compute the *partial sum* of the values
 149 $w(\hat{s}, \hat{s}')$ on the edges of an arbitrary node-to-root path in T ;
- 150 (2) We have to decide how to encode each $w(\hat{s}, \hat{s}')$ on the edges of T . This may affect the
 151 computation of a branching that minimizes such cost;
- 152 (3) We do not know to which k -mer (node of G) each node of T corresponds to.

153 Issues (1) and (2) can be solved with a compressed data structure for answering *partial*
 154 *sums on trees*. In Section 3.1 we first discuss our new compressed solution to partial sums
 155 on the special case of arrays, which may be of independent interest. Then, in Section 3.2
 156 we generalize this solution to partial sums on trees. The solution will introduce a cost in
 157 bits for encoding each weight $w(\hat{s}, \hat{s}')$. The branching $T = (V, E)$ will then be chosen so
 158 as to minimize the sum of the costs on its edges. Issue (3) requires us to keep a mapping
 159 between the nodes of G and T . In Section 3.3 we show how to solve this problem with a
 160 simplified (and practical) variant of the structure of Section 3.2. Finally, in Section 4 we
 161 present experimental results on real DNA datasets, comparing the space usage of our data
 162 structure with the state of the art.

163 3.1 Partial Sums on Arrays

164 Given a (static) sequence $s[1, n]$ of **positive** integers such that $u = \sum_{i=1}^n s[i]$, we would
 165 like to design a data structure to support partial sum operations, i.e., to answer queries
 166 $\text{sum}(i) = \sum_{j=1}^i s[j]$, for any $i \in [n]$. Note that there are succinct data structures that are
 167 able to support **sum** queries in constant time [10, 11, 32] within the information-theoretic
 168 lower bound of $\lceil \log \binom{u}{n} \rceil = n \log \frac{u}{n} + O(n)$ bits. These classic results are summarized in the
 169 following theorem.

170 ► **Theorem 1** ([10, 11, 32]). *Given a (static) sequence $s[1, n]$ of positive integers such that*
 171 *$u = \sum_{i=1}^n s[i]$, there exists a data structure that can answer any **sum** query in $O(1)$ time*
 172 *using $n \log(u/n) + O(n)$ bits of space.*

² In this paper *spanning branching* stands for *spanning forest of arborescences*, that is, a collection of disjoint directed trees spanning the de Bruijn graph. We note that a spanning forest of *undirected* trees would work as well; however, as we discuss in Section 3.3 this introduces some technical complications.

173 This space bound is always at most the space required by writing down all partial sums
 174 explicitly, i.e., $n\lceil\log u\rceil$ bits. However, it is possible to get a better space bound in terms
 175 of a well-known data-aware measure called *gap* measure [15], where the gap measure of the
 176 sequence s is defined as $\mathbf{gap}(s) = \sum_{i=1}^n \lceil\log(s[i] + 1)\rceil$ bits. Notice that in general it is not
 177 possible to represent each $s(i)$ with $\lceil\log(s(i) + 1)\rceil$ bits, and thus data structures must incur
 178 a space overhead on top of the gap measure. Note that the gap measure is always at most
 179 the information-theoretic lower bound of $\lceil\log\binom{u}{n}\rceil$ bits: this is because the gap measure is
 180 maximised when $s[i] = u/n$ for every i . However, in practice the gap measure could be much
 181 smaller than the information-theoretic lower bound.

182 Gupta *et al.* [15] designed a data structure that is able to answer each **sum** query
 183 in $O(\log\log u)$ time and uses $\mathbf{gap}(s) + O(n\log\log(u/n) + n\log(u/n)/\log n)$ bits of space.
 184 Delpratt *et al.* [7] defined the delta measure $\Delta(s) = \sum_{i=1}^n |\delta(s[i])|$, where $|\delta(x)|$, for $x \geq 1$, is
 185 the size in bits of encoding the number x with Elias' δ coding. They present a data structure
 186 that answers each **sum** query in $O(\log\log u)$ time using $\Delta(s) + o(n)$ bits of space. Since for
 187 any x we have $\delta(x) = \log x + 2\log\log x + O(1)$ bits, $\Delta(s) \leq \mathbf{gap}(s) + 2n\log\log(u/n) + O(n)$
 188 bits [7]. This improves over Gupta *et al.* due to the lower order term in the space bound.

189 In the remainder of this section we prove the following theorem. Even if our result on
 190 de Bruijn graphs requires partial sums on trees instead of arrays, we present this result on
 191 arrays here mainly for two reasons. First, it will be used to solve partial sums on trees in
 192 Section 3.1. Second, we believe it can be of independent interest since it provides different
 193 query-time/space trade-offs for partial sums on arrays.

194 ► **Theorem 2.** *Given a (static) sequence $s[1, n]$ of positive integers such that $u = \sum_{i=1}^n s[i]$,
 195 there exists a data structure that answer **sum** queries in $O(\frac{1}{\epsilon})$ time using $(1 + \epsilon)\mathbf{gap}(s) +$
 196 $O(n\log\log(u/n) + u^{\frac{1}{c}})$ bits of space, for any $0 < \epsilon \leq 1$ and any constant $c \geq 1$.*

197 We first present a solution using standard ideas and tools. This solution will be then
 198 refined to obtain the claimed result. We start by concatenating in a vector D the Elias'
 199 δ encoding of each value of s . Elias' γ and δ codes [10] are two standard encodings for
 200 positive integers. Notice that binary encoding could also suffice however we believe that the
 201 proof is clearer using the Elias's one. Elias' $\gamma(x)$ of an integer $x > 0$ is obtained by writing
 202 $\lceil\log x\rceil$ in unary, followed by the value $x - 2^{\lceil\log x\rceil}$ written with $\lceil\log x\rceil$ bits. Elias' δ instead
 203 encodes x by writing $\lceil\log x\rceil + 1$ with Elias' γ , followed by the value $x - 2^{\lceil\log x\rceil}$ written with
 204 $\lceil\log x\rceil$ bits. We would like also to be able to compute the starting position in D of the
 205 encoding of any integers of s . This can be easily done by representing the sequence $L[1, n]$
 206 with the data structure of Theorem 1. The entry $L[i]$ equals the length of the encoding of
 207 $s[i]$. This way, the starting position of the encoding of $s[i]$ in D is $\mathbf{sum}(i)$ in L . Note that
 208 D and L suffice to decompress any value of the original sequence s in constant time. The
 209 space required by these two vectors is bounded by $\mathbf{gap}(s) + O(n\log\log(u/n))$ bits. In order
 210 to support **sum** queries, we can partition s into blocks of size $b = \lceil\frac{1}{\epsilon}\rceil$ each. We store in
 211 an array $B = [1, n/b]$ the partial sums in s up to the beginning of each block. This way
 212 we can easily support any $\mathbf{sum}(i)$ in $O(\frac{1}{\epsilon})$ time. We first get from B the sum up to the
 213 block, say j , that contains the i th element of s . Then, we decompress the elements in the
 214 j th blocks one after the other up to the element $s[i]$. This clearly requires $O(1/\epsilon)$ time per
 215 operation. However, the space needed by B is $\epsilon n \log u$ bits. This gives an overall space usage
 216 of $\mathbf{gap}(s) + \epsilon n \log u + O(n\log\log(u/n) + n)$ bits, which is worse than what is claimed in
 217 Theorem 2, since the term $\epsilon n \log u$ may be potentially larger than the term $\epsilon \mathbf{gap}(s)$.

218 To improve the space bound, we use a different partitioning strategy. The goal is to
 219 partition s into variable size blocks such that the cost of encodings the values in each block

(but the last one) is between $\frac{\log u}{\epsilon}$ bits and $\frac{3 \log u}{\epsilon}$ bits. This guarantees that there are at most $\frac{\epsilon \text{gap}(s)}{\log u} + 1$ blocks and, thus, the cost of storing the vector B is at most $\epsilon \text{gap}(s) + \log u$ bits.

We observe that a partition with the above characteristics is always possible. For example we can use the following greedy algorithm. We start with an empty block and we process the sequence from left to right. While processing an element we have two possibilities: we either include the element in the current block or we create a new block with this element. We take the former decision only if the overall cost of encoding the elements in the block is less than $\frac{\log u}{\epsilon}$ bits. As the code of an element is at most $\log u + O(\log \log u)$ bits, the above partitioning strategy gives blocks of size between $\frac{\log u}{\epsilon}$ bits and $\frac{3 \log u}{\epsilon}$ bits as claimed.

As now we have variable-size blocks, we need to store a binary vector V of size n bits to keep track of blocks boundaries: an entry $V[i]$ is 1 if and only if the i th element of s is the first element in a block. We use a data structure to support rank/select operations on V in constant time [32].

A query $\text{sum}(i)$ is answered as follows. We first use rank/select operations on V to compute the block j of position i and its offset p within this block. We now need to decode the first p elements of the j th block and sum them to the value $B[j]$. Unfortunately, there may be $\Theta(\log u)$ elements encoded in a block and, thus, we cannot use a trivial decoding. Instead, we conceptually split each block into subblocks such that each subblock is either (i) formed of elements whose overall encoding sizes is no more than $\frac{1}{2} \log u$ bits, or (ii) a single element. This split into subblocks can be done by an easy variant of the greedy partitioning strategy above. Note that there are $O(\frac{1}{\epsilon})$ subblocks per block. We use a table of size $O(\sqrt{u} \log u)$ bits to precompute the sum of any prefix of any subblock given its encoding. This way, computing the sum of the first p elements in j th block costs $O(\frac{1}{\epsilon})$ time as required. Note that the splitting above can be changed to use $\frac{1}{c} \log u$ bits, for any constant $c \geq 1$ instead $\frac{1}{2} \log u$ bits. This way we need a table of size $O(2^{\frac{1}{c} \log u} \log u) = O(u^{\frac{1}{c}} \log u)$ bits and the query time remains $O(\frac{1}{\epsilon})$. By adjusting infinitesimally the constant c , the table takes $O(u^{\frac{1}{c}})$ bits of space.

We conclude by showing how to adapt any solution for the partial sums problem to the variant of the problem in which the sequence s has both positive and negative integers. Let us define $\text{gap}^{\pm}(s) = \sum_{i=1}^n \lceil \log(|s[i]| + 1) \rceil$

This variant can be easily reduced to two instances of the original problem as follows. We first use a binary vector $S[1, n]$ that records the signs of the values in s , i.e., $S[i] = 1$ if and only if $s < 0$, $S[i] = 0$ otherwise. We use the data structure by Raman *et al.* [32] to support constant time rank/select operation on S using $\log \binom{n}{p} + o(n)$ bits of space, where p is the number of positive integers in s . Then, we create two sequences $s^+[1, p]$ and $s^-[1, n - p]$ that store positive and negative integers of s , respectively. Any partial sum query can now be answered with two partial sums queries on s^+ and s^- . Combining this reduction with Theorem 2 yields the following corollary.

► **Corollary 3.** *Given a (static) sequence $s[1, n]$ of positive and negative integers such that $u = \sum_{i=1}^n |s[i]|$, there exists a data structure that answer each sum operation in $O(\frac{1}{\epsilon})$ time using $(1 + \epsilon) \text{gap}^{\pm}(s) + \log \binom{n}{p} + O(n \log \log(u/n) + u^{\frac{1}{c}})$ bits of space, for any $0 < \epsilon \leq 1$ and any constant $c \geq 1$, where p is the number of positive integers.*

3.2 Partial Sums on Trees

Let $T = (V, E)$ be a rooted weighted tree and let $w(i, j)$ denote the (possibly negative) integer weight of edge $(i, j) \in E$. We would like to build a space-efficient data structure that answers partial sum queries of the form $\text{sum}(v) = \sum_{(i, j) \in \Pi(v)} w(i, j)$ on paths, i.e., that

266 reports the sum of the weights on the path $\Pi(v)$ connecting node v to the root, where v is
 267 represented in pre-order (the solution for post-order is symmetrical and we do not discuss it).

268 Chan *et al.* [5] tackle this problem in a more general framework, where weights belong to
 269 a semigroup of size γ . Let $n = |V|$ be the number of nodes. In this model, they provide a
 270 data structure taking $n \log \gamma + o(n \log \gamma) + 2n$ bits and answering queries in $O(\alpha(n))$ (inverse
 271 Ackermann) time. Note that this space is succinct but does not achieve compression: each
 272 weight $w(i, j)$ is stored in $\log \gamma$ bits, independently of its magnitude. Typically, one would like
 273 a data structure using much less space: ideally, close to $\text{gap}^\pm(T) = \sum_{(i,j) \in E} \lceil \log(|w(i, j)| + 1) \rceil$
 274 bits of space.

275 A folklore approach can reduce partial sums on trees to partial sums on arrays. The
 276 idea consists of linearizing the weights according to a Euler tour of the tree. Recall that
 277 tree edges are directed towards the root and note that the Euler tour visits each tree edge
 278 (i, j) twice: the first time in the direction (j, i) and the second time in the direction (i, j) .
 279 If $(i, j) \in E$, we assign to the reverse edge (j, i) weight $w(j, i) = -w(i, j)$. We initialize
 280 an empty sequence W and, for each edge (i, j) visited along the tour (that is, in the direction
 281 $i \rightarrow j$), append $-w(i, j)$ at the end of W . Then, it holds that $\text{sum}(v) = \sum_{j=1}^t W[j]$, where
 282 t is the index corresponding to the first time we see node v in the Euler tour (i.e., $W[t]$
 283 contains the weight $w(\pi(v), v)$, where $\pi(v)$ is the parent of v). This equality holds because,
 284 in $W[1], \dots, W[t]$, all values that correspond to edges not belonging to the path from v to
 285 the root appear two times with opposite signs and thus they cancel out.

286 We can use the data structure of Corollary 3 on the sequence W to have $O(\frac{1}{\epsilon})$ query
 287 time. The space is $(2 + \epsilon)\text{gap}^\pm(s) + O(n \log \log(u/n) + u^{\frac{1}{c}})$ bits because every gap occurs
 288 twice in W .

289 Another folklore approach uses the Heavy-Light decomposition of the dynamic trees
 290 of Sleator and Tarjan [34]. The idea is to split the tree into heavy paths and use a data
 291 structure for prefix sums on each of these heavy paths. By virtues of the decomposition,
 292 any root-to-a-node path crosses $O(\log n)$ heavy paths and, thus, a query can be answered
 293 with $O(\log n)$ prefix sums on arrays. This solution has a better space bound as every gap is
 294 represented exactly once, but its query time is logarithmic.

295 In the remainder of this section we prove the following theorem, which improves over the
 296 two folklore solutions above.

297 ► **Theorem 4.** *Given a (static) rooted tree $T = (V, E)$ with positive integer weights $w(i, j)$
 298 associated with each edge $(i, j) \in E$ such that $u = \sum_{(i,j) \in E} w(i, j)$, there exists a data
 299 structure that answers sum queries in $O(\frac{1}{\epsilon})$ time using $(1 + \epsilon)\text{gap}(T) + O(n \log \log(u/n) + u^{\frac{1}{c}})$
 300 bits of space, for any $0 < \epsilon \leq 1$ and any constant $c \geq 1$.*

301 The following corollary generalizes the solution to trees with positive and negative weights.

302 ► **Corollary 5.** *Given a (static) rooted tree $T = (V, E)$ with positive and negative integer
 303 weights $w(i, j)$ associated with each edge $(i, j) \in E$ such that $u = \sum_{(i,j) \in E} |w(i, j)|$, there
 304 exists a data structure that solves queries sum in $O(\frac{1}{\epsilon})$ time using $(1 + \epsilon)\text{gap}^\pm(T) + \log \binom{n}{p} +$
 305 $O(n \log \log(u/n) + u^{\frac{1}{c}})$ bits of space, for any $0 < \epsilon \leq 1$ and any constant $c \geq 1$, where p is
 306 the number of positive integers.*

307 These results are obtained by adapting to trees the solutions of Theorem 2 and Corollary 3,
 308 respectively.

309 The high level ideas behind our approach are as follows. We first partition the tree into
 310 subtrees. Subtrees are either disjoint or intersect only at their common root. Then, we
 311 store the sums on the paths from the tree's root to each subtree root. The sum $\text{sum}(v)$ is

312 computed by taking the sum up to the root of v 's subtree and summing up the sum of the
 313 path within the subtree. The sum within each subtree is computed by using an approach
 314 similar to the one in the previous subsection.

315 The partitioning of the tree is crucial for the space and time efficiency of the solution.
 316 For this aim we use a weighted variant of the tree covering procedure described by Geary *et al.*
 317 *al.* [14, Sec. 2.1] to decompose T into subtrees. Given any parameter $M > 2$, the original
 318 tree covering procedure as described in [14, Sec. 2.1] is used to decompose T into $\Theta(n/M)$
 319 sub-trees containing $O(M)$ nodes each. Two subtrees are either disjoint or intersect only at
 320 their common root. The covering is built by visiting the tree in post-order and by grouping
 321 nodes into components. The procedure lets the current component grow until its size falls in
 322 the range $[M, 3M - 4]$. When a component reaches the required size, a new empty component
 323 is created. This greedy procedure guarantees that every component (except for the one
 324 containing the root) has between M and $3M - 4$ nodes. See Geary *et al.* [14, Sec. 2.1] for
 325 more details.

326 In our solution we use a simple variant of this covering procedure. We fix M to be $\frac{\log u}{\epsilon}$
 327 and we use the greedy procedure above with the only difference that every node i has an
 328 encoding cost which equals the size of encoding $w(\pi(i), i)$, where $\pi(i)$ is the parent of i . We
 329 build our components to have a cost bounded by M . This way, we can control the size of the
 330 encoding of each subtree, which is between $\frac{\log u}{\epsilon}$ bits and $\frac{3 \log u}{\epsilon} + 2 \log u$ bits. We now need
 331 to δ -encode the weights w of a subtree one after the other by following an ordering (specified
 332 later) which is suitable for solving queries efficiently. Our goal is, given a node v represented
 333 in pre-order, to compute the sum from tree's root to node v . This is solved by summing up:
 334 i) the sum from the tree's root to the root of the subtree containing v ; ii) the sum of the
 335 weights on the path from the subtree's root to node v . Geary *et al.* [14, Sec. 4.3] show how
 336 to find the pre-order number of the root of the subtree containing v in constant time and
 337 $o(n)$ bits of space. This index can be used to access a vector containing the sum mentioned
 338 at point (i). Point (ii) is computed in $O(\frac{1}{\epsilon})$ time in a way similar to what we have done for
 339 arrays. We split the representation of a subtree into $O(\frac{1}{\epsilon})$ subblocks and use tables storing
 340 precomputed answers to all the possible subblocks representations. However, the solution
 341 here is more involved than the one we used for arrays because we need to compute the sum
 342 only of some of the elements of a subblock, i.e., those elements that belong to the query
 343 path, and exclude the other elements. This can be done by using a mask that marks the
 344 elements belonging to the query path. Given the representation of a subblock and a mask,
 345 a precomputed table returns the sum of the marked elements only. The main issue is now
 346 to compute the required mask. This can be done by splitting the subtree into subsubtrees
 347 whose encoding takes between $\frac{1}{12} \log u$ bits and $\frac{1}{4} \log u$ bits. This is done by using once again
 348 the covering procedure described above. We δ -encode the elements in these subsubtrees
 349 by visiting nodes in BFS order. We also need to store the topology of each subsubtree.
 350 To do that, we can use a balanced sequence representation of each subsubtree so that we
 351 need $2n + o(n)$ bits overall. Observe that each subsubtree has at most $\frac{1}{4} \log u$ nodes. Thus,
 352 its balanced parentheses representation always takes at most $\frac{1}{2} \log u$ bits and we can use a
 353 table of size $O(u^{\frac{1}{2}} \log^2 u)$ bits that, given a subsubtree topology ($\frac{1}{2} \log u$ bits) and a node
 354 represented as a pre-order position in the subsubtree topology ($O(\log \log u)$ bits), returns the
 355 required mask, i.e., a mask of $\frac{1}{4} \log u$ bits that marks only nodes on the path from the root
 356 to the node. Another table of size $O(u^{\frac{1}{2}} \log u)$ bits is used to, given any possible combination
 357 of the delta-encoded weights of a subsubtree (at most $\frac{1}{4} \log u$ bits) and any possible mask
 358 (at most $\frac{1}{4} \log u$ bits), return the sum of the marked elements only. Note that to perform
 359 the above operations we need, given a pre-order node v , to obtain (a) the packed balanced

360 parentheses sequence representation of the subtree containing v and (b) the pre-order
 361 position of v within its subtree. Information (a) can be obtained as done above by using
 362 the procedure described by Geary et al. [14, Sec. 4.3]: we locate the pre-order number of
 363 the root of the subtree and use it as index in an array containing the packed balanced
 364 sequence representations of the subtrees ($O(n)$ bits of space). In the same section, Geary
 365 et al. [14, Sec. 4.3] show also how to obtain information (b) in constant time and $o(n)$ bits
 366 of additional space. Finally note that, as done in the previous section, we can replace the
 367 size $\frac{1}{4} \log u$ of the subtrees by $\frac{1}{c} \log u$ for any constant $c \geq 1$ and obtain the claimed space
 368 bound.

369 3.3 Adding the dBg Topology

370 In this section we discuss how to combine a simplified (and practical) variant of the data
 371 structure of Corollary 5 with the BOSS de Bruijn graph representation of Bowe et al. [3].
 372 Our final data structure represents a weighted de Bruijn graph in compressed space and
 373 supports computing the abundance of any given input k -mer. In Section 4 we present
 374 experimental results based on our data structure. Our implementation is available at
 375 <https://github.com/nicolaprezza/cw-dBg/>.

376 Let n and m be the number of nodes and edges of the input de Bruijn graph G , respectively.
 377 First of all, we compute a branching T of G minimizing measure $\text{gap}^\pm(T)$. This can be
 378 achieved using Gabow et al.'s optimized version [13] of Edmonds' algorithm [9], running in
 379 $O(m + n \log n)$ time.

380 The second step is to observe that, by definition of branching, its topology is embedded in
 381 the topology of the de Bruijn graph. Consider the list $\hat{s}_1, \dots, \hat{s}_n$ of the n nodes of the de Bruijn
 382 graph sorted by the co-lexicographic order of their corresponding k -mers. Note that this is
 383 precisely the order in which nodes are stored in the BOSS representation [3]. Let $e_i^1, \dots, e_i^{t_i}$
 384 be the t_i incoming edges of node \hat{s}_i , and consider the list $e_1, \dots, e_m = \langle e_i^1, \dots, e_i^{t_i} \rangle_{i=1, \dots, n}$ of
 385 all the m edges in the graph, sorted by their target node. We keep one bitvector $\text{IN_B}[1, m]$
 386 marking the edges, in the above order, that are included in the branching. It is clear that
 387 the topology of the de Bruijn graph, combined with bitvector IN_B , fully specifies the
 388 topology of the branching. Importantly, we observe that de Bruijn graphs are usually very
 389 sparse, i.e., $m \approx n$. This implies that IN_B is composed mostly of bits equal to 1, thus
 390 its zero-order entropy $\lceil \log \binom{m}{n} \rceil$ is likely to be very small. We thus store IN_B with the
 391 zero-order compressed representation of Raman et al. [32], supporting constant-time rank
 392 and select queries.

393 The third step is to build a simplified version of the structure of Corollary 5 over each
 394 arborescence of the branching. We do this as follows. Given a parameter ρ (the *sample*
 395 *rate*), $1 \leq \rho \leq n$, we use the tree covering procedure of Geary et al. [14, Sec. 2.1] to
 396 decompose each arborescence into $\Theta(n'/\rho)$ subtrees of $O(\rho)$ nodes each, where n' is the
 397 number of nodes of the arborescence. We furthermore explicitly store the abundances of the
 398 subtree's roots in a vector, sorting them by the order in which nodes appear in the BOSS
 399 data structure (that is, by the co-lexicographic order of their corresponding k -mers). The
 400 sampled abundances take $O((\frac{n}{\rho} + t) \log u)$ bits, where t is the number of arborescences. An
 401 additional zero-order compressed bitvector marks sampled nodes. Using the representation
 402 of Raman et al. [32], this bitvector takes $O((\frac{n}{\rho} + t) \log \rho) + o(n)$ bits and supports access,
 403 rank, and select operations in constant time.

We store the weight of each edge of T (that is, the difference between the abundances
 of its endpoint k -mers) using a version of Elias' gamma encoding supporting fast random
 access queries. First, we convert each (possibly negative) weight w into a positive integer

23:10 Compressed Weighted de Bruijn Graphs

$w' \geq 1$ using the formula

$$w' = \begin{cases} -2w, & \text{if } w < 0 \\ 2w + 1, & \text{otherwise} \end{cases}$$

404 In practice, using the above conversion, rather than storing explicitly the sign of w in
 405 a separate bitvector, has an important practical advantage, as the minimum branching
 406 algorithm compresses also the sign of the weights. Indeed, we observed experimentally that
 407 this choice compresses each sign to about 0.3 bits on average (rather than 1 bit per sign
 408 required when storing them in a plain bitvector). The random access gamma-compressed
 409 weights are implemented as follows. We concatenate the binary representation of each w'
 410 (devoid of its most significant bit) in an uncompressed bitvector X , and its length in unary
 411 in a bitvector Y (for example: $w' = 27$ would be encoded as 1011 in X and 10000 in Y). We
 412 compress Y using the representation of Raman et al. [32]. The i -th weight can be extracted
 413 from this representation in $O(1)$ time with one select operation on Y and one packed access
 414 operation on X . Since we compress Y , it is not hard to see that our representation takes
 415 at most $\text{gap}^\pm(T) + O(n \log \log(u/n))$ bits of space. Finally we mention that, in order to
 416 achieve further compression, we used the compressed bitvector of Raman et al. [32] also to
 417 implement the components of the BOSS representation [3].

418 Let s be a k -mer. To retrieve its abundance, we do the following:

- 419 1. We use the de Bruijn graph representation to retrieve the node \hat{s} (in the BOSS repre-
 420 sentation, a position in the Burrows-Wheeler transform of the graph) corresponding to
 421 s . This step takes $O(k)$ time via the backward search algorithm (see [3] for full details)
 422 since we assume constant-sized DNA alphabet $\Sigma = \{A, C, G, T\}$.
- 423 2. Starting from \hat{s} , we move upward towards the root of the tree containing \hat{s} in the branching,
 424 stopping as soon as a sampled node \hat{s}' is found (that is, a node whose abundance has been
 425 stored explicitly). Each move in the tree is implemented with a constant-time application
 426 of the FL function [3]. Along the walk, we sum the abundance of the sampled node \hat{s}'
 427 to the weights of the edges on the path connecting \hat{s}' with \hat{s} . Overall, this step takes $O(\rho)$
 428 time.

429 Observe that Step 1 can be avoided if we are already given the representation of a node
 430 in the de Bruijn graph (for example, if we are navigating it). Now we can also explain why
 431 we use a branching instead of a minimum spanning undirected forest. With a branching,
 432 the parent of node \hat{s} (in its corresponding arborescence) is always one of its incoming edges
 433 (precisely, the one marked in bitvector `IN_B`). With an undirected spanning forest, instead, the
 434 parent of node \hat{s} could be one of its outgoing edges; this would force us to use an additional
 435 bitvector, increasing the overall space and slowing down operations.

436 To sum up, our implementation offers the following trade-offs. Let G be a de Bruijn
 437 graph with m edges and n nodes, $w(\hat{i}, \hat{j}) = c(\hat{i}) - c(\hat{j})$ be the weight associated with each
 438 edge (\hat{i}, \hat{j}) of G , where $c(\hat{i})$ is the abundance of node \hat{i} , $T = (V, E)$ be a branching of G
 439 with t connected components minimizing measure $\text{gap}^\pm(T) = \sum_{(\hat{i}, \hat{j}) \in E} [\log(|w(\hat{i}, \hat{j})| + 1)]$,
 440 $u = \sum_{(\hat{i}, \hat{j}) \in E} |w(\hat{i}, \hat{j})|$, and $1 \leq \rho \leq n$ be the user-defined sample rate. Our data structure
 441 uses $\text{gap}^\pm(T) + O\left(n \log \log(u/n) + \left(\frac{n}{\rho} + t\right) \log u\right) + \log \binom{m}{n}$ bits on top of the compressed
 442 BOSS representation [3] and allows us to retrieve:

- 443 (1) the abundance of a given k -mer $s \in \Sigma^k$ in $O(k + \rho)$ time, and
- 444 (2) the abundance of a given node in the de Bruijn graph (represented as a position in BOSS)
 445 in $O(\rho)$ time.

446 Note that, in practice, ρ is usually chosen to be $\rho \gg k$ (in our experiments, $\rho = 64$ and
447 $k = 28$) so the two query times are not expected to differ much in practical applications.

448 Our experiments (see next section) highlight that, in practice, the number t of arbores-
449 cences is negligible compared to the size of the graph.

450 4 Experiments

451 In order to get a feeling of its practical value, we implemented a first preliminary version of
452 our compressed representation of weighted de Bruijn graphs. In the following, we refer to
453 this implementatio as `cw-dBg`; its code is available at [https://github.com/nicolaprezza/
454 cw-dBg](https://github.com/nicolaprezza/cw-dBg).

455 **Tools.** We ran `cw-dBg` on four datasets and compared the results obtained against the
456 state-of-the-art algorithms for this problem: `Squeakr` [29], `deBGR` [27], and the very recent tool
457 `fress` [33] which appeared online in November 2020. Both `Squeakr` and `deBGR` are based on
458 the so-called *counting quotient filter (CQF)* data structure [28]. `Squeakr` supports two modes:
459 approximate and exact. When run in exact mode, the k -mers are inserted in the CQF using
460 an invertible $2k$ -bit hash function. In approximate mode, `Squeakr` uses a p -bit hash function,
461 with $p \leq 2k$, and thus can have false positives with an error rate depending on the chosen
462 p . `deBGR` is based on `Squeakr` and builds an approximate data structure with smaller error
463 rate by increasing the space complexity of approximate `Squeakr` by just 18% – 20%. Finally,
464 `fress` [33] implements an approximate data structure called *Set-Min sketch*, inspired by the
465 *Count-Min sketch* data structure [6], that takes advantage of the power-law distribution of
466 the k -mer counts to reduce both the error rate of the returned results and the space usage
467 (by an entropy-compression mechanism). Similarly to Count-Min sketch, Set-Min sketch uses
468 several hash functions to map a given k -mer to sketches of its abundance; at query time,
469 the abundance of the k -mer is computed by combining the retrieved sketches. As shown by
470 the authors, Set-Min sketch is more space-efficient than minimal perfect hash functions and
471 provides better error guarantees compared to equally-dimensional Count-Min sketches.

472 **Datasets.** We selected datasets to cover applications that could be as widely different as
473 possible. In particular, we tested the above tools against the following datasets:

- 474 ■ `IAVs Inf. 1 (1.2Gbp)` and `IAVs Inf. 2 (578Mbp)` are samples from the dataset presented in
475 [1] of human lung cells infected by Influenza A viruses (IAVs) and correspond to a sample
476 of 10 millions reads and a sample of 4,814,148 reads from chromosome 2, respectively.
- 477 ■ `E. coli (2.3Gbp)` consists of 22,720,100 Illumina reads of *E. Coli* K-12 strain MG1655
478 (available in the ENA repository under the following study: PRJEB2323, [https://www.
479 ebi.ac.uk/ena/browser/view/ERR022075](https://www.ebi.ac.uk/ena/browser/view/ERR022075)).
- 480 ■ `Human (6.5Gbp)` consists of 63,917,134 Illumina reads of human RNA (available in the
481 ENA repository under the following study: PRJNA609878, [https://www.ebi.ac.uk/
482 ena/browser/view/SRX7829390](https://www.ebi.ac.uk/ena/browser/view/SRX7829390)).

483 **Experimental settings.** The experiments were carried out on a single core of an 8-core
484 Intel i9-9900K server with 64 GB of RAM and running Linux 5.4.0, 64 bits. Our code is
485 written in C++ and compiled with `gcc 9.3.0`.

486 As `deBGR` assumes k -mers of size 28 (which cannot be changed), for consistency we also
487 used this value in all our experiments for all tools. Table 1 reports some characteristics of
488 the datasets. Notice that, as it was previously mentioned, the number of arborescences is

23:12 Compressed Weighted de Bruijn Graphs

489 very small, when compared to the size of the graph, and thus negligible in practice. We used
 490 sample rate $\rho = 64$ for cw-dBg.

Dataset	# bases	# k -mers	# edges	# Arb.	Avg ab.	Max ab.	# ab.
IAVs Inf. 1	1,200,000,000	57,629,309	65,880,010	259	16.14	607,028	13,452
IAVs Inf. 2	577,697,760	34,914,869	41,295,094	100	12.82	55,110	9,007
E. coli	2,317,450,200	87,414,957	281,845,871	1	19.49	8,520,260	2,657
Human	6,455,630,534	362,818,017	463,413,958	1226	13.04	11,054,076	21,472

■ **Table 1** Characteristics of the de Bruijn graph resulted from each of the dataset for k -mers of size 28. The columns correspond to: number of bases and number of distinct k -mers in the dataset, number of edges and number of arborescences in the de Bruijn graph, average and max abundance of the k -mers and the total number of distinct abundances in the dataset (which influences `fress`’ space).

491 deBGR and Squeakr (the latter both in its approximate and exact version) were run using
 492 only 1 thread. When running deBGR, the user needs to specify two parameters related to the
 493 log of the number of slots in the counting quotient filter. The *CQF size* argument is estimated
 494 using the `lognumslots.sh` script provided by the authors and corresponds approximately to
 495 the log of the number of k -mers. The *exact CQF size* argument is estimated as 20% less
 496 then the *CQF size* (personal communication with the authors of deBGR).

497 Finally, `fress` is an approximate tool and allows to specify a parameter related to the
 498 error rate: in our experiments, we ran it using an error rate of $\epsilon = 0.01$, as used in `fress`’
 499 paper [33]. Notice that the error rate estimates the expected number of collision which could
 500 possibly lead to output a wrong abundance. An error rate of 0.01 means that in N queries,
 501 the number of expected collisions in the hash table is ϵN .

502 **Experimental results.** In Table 2 we present the detailed space required by our represen-
 503 tation for each of the datasets considered, while Table 3 reports the space usage of all the
 504 tools considered.

505 The bit/ k -mer of each component in Tables 2 and 3 were computed using the number
 506 of distinct k -mers over the alphabet $\{A, C, G, T\}$ present in the original dataset. We note
 507 that this is not the same calculation performed by Bowe et al. [3] when estimating the size
 508 of their BOSS representation: in that case, the authors divide by the number of *edges* in
 509 the de Bruijn graph when also including dummy (that is, left-padded) k -mers, obtaining
 510 roughly 4 bits per edge. In fact, the two values (number of edges and number of k -mers) are
 511 close to each other only when the number of dummy k -mers is small (which is not always
 512 the case, as explained below). We believe that dividing the total bit-size by the number
 513 of distinct k -mers appearing in the dataset is more general, since it applies to any k -mer
 514 counting data structure like the ones shown in Table 3. Notice that for this reason in Table
 515 2 the space reported for the BOSS representation of the E. coli dataset is larger than the
 516 typical 4 bit/edge reported by the authors of BOSS [3] as the number of dummy k -mers for
 517 this dataset is significantly large (approximately equal to the number of distinct k -mers).

518 The results in Table 2 and Table 3 show that our representation uses at most 4.75
 519 bit/ k -mer (and as little as 1.42 bit/ k -mer) to encode the abundances. Even when adding the
 520 BOSS representation (required to map k -mers to their corresponding abundances), we obtain
 521 at most 12.64 bit/ k -mer (and as little as 4.14 bit/ k -mer) while the exact version of Squeakr
 522 uses 83.26 bit/ k -mer in the best scenario. In almost all our experiments cw-dBg required
 523 substantially much less space, an by several orders of magnitude, than all the other tools.
 524 The only exception occurred in the E. coli dataset, where `fress` performed slightly better than

Dataset	Final size (MB)	BOSS (bit/ k -mer)	Compressed ab. (bit/ k -mer)			Total size (bit/ k -mer)	Query time (μ s/query)
			weights	branching	total		
IAVs Inf. 1	29	2.73	0.95	0.48	1.43	4.15	47.56
IAVs Inf. 2	20	2.87	1.19	0.53	1.72	4.59	44.24
E. coli	132	7.88	2.93	1.82	4.75	12.64	62.34
Human	223	3.14	1.40	0.61	2.01	5.14	67.42

■ **Table 2** The space required for each dataset by our representation of the weighted de Bruijn graph. The columns correspond to: the size of the file containing the final representation, the space required by the BOSS representation of the de Bruijn graph; the space required by the compressed abundances divided in the space required by the delta-compressed weights (stored on the edges of the branching) and the one for the branching topology and the sampled abundances on the root of each subtree; the total space required by our representation and the average query time in microseconds for retrieving the abundance of a given k -mer (represented as a string). Space is measured in bits per distinct k -mer (see Table 1) and we used sample rate $\rho = 64$.

	IAVs Inf. 1		IAVs Inf. 2		E. coli		Human	
	MB	bit/ k -mer	MB	bit/ k -mer	MB	bit/ k -mer	MB	bit/ k -mer
cw-dBg	29	4.15	20	4.59	132	12.64	223	5.14
Squeakr approx.	325	47.19	179	42.80	325	31.11	1126	24.05
Squeakr exact	965	140.40	499	119.75	965	92.57	3686	83.26
deBGR	912	132.46	376	89.55	912	87.32	5529	119.00
fress	733	106.56	135	32.35	115	10.95	733	16.90

■ **Table 3** For each dataset we report the space required for the representation of the weighted de Bruijn graph by each of the considered tools.

525 cw-dBg by using 10.95 bit/ k -mer (versus 12.64 bit/ k -mer of cw-dBg). We stress out, however,
526 that **fress** achieves this slight space saving at the cost of returning a wrong abundance 1% of
527 the times in the worst case. Indeed, if its error rate is lowered from 0.01 to 0.005, **fress** uses
528 14.6 bit/ k -mer on this dataset. The efficiency of **fress** on this dataset is explained by the
529 fact that the number of distinct k -mer abundances is much smaller than in the other three
530 datasets (see Table 1). This is clearly the case where we expect **fress** to work well, as it fully
531 takes advantage of the power-law distribution of the abundances by entropy-compressing
532 their values.

533 Finally, we estimated the average query time of cw-dBg, by querying a set of 5,000
534 randomly chosen 28-mers. The largest measured query time for cw-dBg on our machine
535 was 67.42 μ s/query on the **Human** dataset (see Table 2 for the query times on all datasets).
536 As expected, hash-based data structures have better query time: on the same **Human**
537 dataset, **Squeakr** answers queries in average 0.19 μ s/query using 1 thread. However, as we
538 discussed before, hash-based structures have a much higher space usage, in addition to being
539 approximate. The exact variant of **Squeakr** uses orders of magnitude more space than cw-dBg.

540 5 Conclusions and open problems

541 We propose a new compressed representation for weighted de Bruijn graphs based on the idea
542 of delta-encoding the variations of k -mer abundances on a spanning branching of the graph.
543 As a by-product of independent interest, we exhibit efficient compressed data structures for

544 answering partial sums on edge-weighted trees.

545 We show both theoretically and experimentally that our approach uses significantly less
 546 memory than the one used by the state-of-the-art exact representations. Our de Bruijn graph
 547 representation is general, in other words it is not restricted by the application (e.g., variation
 548 finding or RNA-seq), and can be used as part of any algorithm that represents NGS data
 549 with de Bruijn graphs. Future extensions will include implementing a strategy similar to
 550 the one used by `fress` in order to take advantage of the small number of distinct abundance
 551 observed in practice. The basic idea is to use a table storing all the distinct abundances
 552 sorted and, for each node, encode its abundance's offset in the table with our data structure
 553 (that is, storing the deltas of these offsets rather than of the original abundances). We also
 554 plan to integrate our structure in a usable bioinformatics tool.

555 — References —

- 556 **1** Usama Ashraf, Clara Benoit-Pilven, Vincent Navratil, Cécile Ligneau, Guillaume Fournier,
 557 Sandie Munier, Odile Sismeiro, Jean-Yves Coppée, Vincent Lacroix, and Nadia Naffakh.
 558 Influenza virus infection induces widespread alterations of host cell splicing. *NAR Genomics
 559 and Bioinformatics*, 2(4), 11 2020.
- 560 **2** Anton Bankevich, Sergey Nurk, Dmitry Antipov, Alexey A. Gurevich, Mikhail Dvorkin,
 561 Alexander S. Kulikov, Valery M. Lesin, Sergey I. Nikolenko, Son Pham, Andrey D. Prjibelski,
 562 Alexey V. Pyshkin, Alexander V. Sirotkin, Nikolay Vyahhi, Glenn Tesler, Max A. Alekseyev,
 563 and Pavel A. Pevzner. Spades: A new genome assembly algorithm and its applications
 564 to single-cell sequencing. *Journal of Computational Biology*, 19(5):455–477, 2012. PMID:
 565 22506599. doi:10.1089/cmb.2012.0021.
- 566 **3** Alexander Bowe, Taku Onodera, Kunihiko Sadakane, and Tetsuo Shibuya. Succinct de Bruijn
 567 Graphs. In Ben Raphael and Jijun Tang, editors, *Algorithms in Bioinformatics*, pages 225–235,
 568 Berlin, Heidelberg, 2012. Springer Berlin Heidelberg.
- 569 **4** Michael Burrows and David J Wheeler. A block-sorting lossless data compression algorithm.
 570 Technical Report 124, Digital Equipment Corporation, 1994.
- 571 **5** Timothy M. Chan, Meng He, J. Ian Munro, and Gelin Zhou. Succinct indices for path minimum,
 572 with applications. *Algorithmica*, 78(2):453–491, June 2017. doi:10.1007/s00453-016-0170-7.
- 573 **6** Graham Cormode and S. Muthukrishnan. Summarizing and mining skewed data streams. pages
 574 44–55, 2005. 5th SIAM International Conference on Data Mining, SDM 2005 ; Conference
 575 date: 21-04-2005 Through 23-04-2005.
- 576 **7** O’Neil Delpratt, Naila Rahman, and Rajeev Raman. Compressed prefix sums. In *Proceedings
 577 of the 33rd Conference on Current Trends in Theory and Practice of Computer Science
 578 (SOFSEM)*, volume 4362 of *Lecture Notes in Computer Science*, pages 235–247. Springer, 2007.
 579 doi:10.1007/978-3-540-69507-3_19.
- 580 **8** Sebastian Deorowicz, Marek Kokot, Szymon Grabowski, and Agnieszka Debudaj-Grabysz.
 581 KMC 2: fast and resource-frugal k-mer counting. *Bioinformatics*, 31(10):1569–1576,
 582 01 2015. arXiv:[https://academic.oup.com/bioinformatics/article-pdf/31/10/1569/
 583 17085507/btv022.pdf](https://academic.oup.com/bioinformatics/article-pdf/31/10/1569/17085507/btv022.pdf), doi:10.1093/bioinformatics/btv022.
- 584 **9** Jack Edmonds. Optimum branchings. *Journal of Research of the national Bureau of Standards
 585 Section B*, 71(4):233–240, 1967.
- 586 **10** Peter Elias. Efficient storage and retrieval by content and address of static files. *J. ACM*,
 587 21(2):246–260, April 1974. doi:10.1145/321812.321820.
- 588 **11** Robert Mario Fano. On the number of bits required to implement an associative memory.
 589 memorandum 61. *Computer Structures Group, Project MAC, MIT, Cambridge, Mass., nd*,
 590 1971.
- 591 **12** Paolo Ferragina and Giovanni Manzini. Indexing compressed text. *J. ACM*, 52(4):552–581,
 592 2005. doi:10.1145/1082036.1082039.

- 593 13 Harold N Gabow, Zvi Galil, Thomas Spencer, and Robert E Tarjan. Efficient algorithms
594 for finding minimum spanning trees in undirected and directed graphs. *Combinatorica*,
595 6(2):109–122, 1986.
- 596 14 Richard F Geary, Rajeev Raman, and Venkatesh Raman. Succinct ordinal trees with level-
597 ancestor queries. *ACM Transactions on Algorithms (TALG)*, 2(4):510–534, 2006.
- 598 15 Ankur Gupta, Wing-Kai Hon, Rahul Shah, and Jeffrey Scott Vitter. Compressed data
599 structures: Dictionaries and data-aware measures. *Theor. Comput. Sci.*, 387(3):313–331, 2007.
600 doi:10.1016/j.tcs.2007.07.042.
- 601 16 David L. Hyten, Steven B. Cannon, Qijian Song, Nathan Weeks, Edward W. Fickus, Randy C.
602 Shoemaker, James E. Specht, Andrew D. Farmer, Gregory D. May, and Perry B. Cregan.
603 High-throughput snp discovery through deep resequencing of a reduced representation library
604 to anchor and orient scaffolds in the soybean whole genome sequence. *BMC Genomics*, 11(1):38,
605 2010. doi:10.1186/1471-2164-11-38.
- 606 17 Katerina Kechris, Yee Hwa Yang, and Ru-Fang Yeh. Prediction of alternatively skipped
607 exons and splicing enhancers from exon junction arrays. *BMC Genomics*, 9(1):551, 2008.
608 doi:10.1186/1471-2164-9-551.
- 609 18 Ruiqiang Li, Yingrui Li, Xiaodong Fang, Huanming Yang, Jian Wang, Karsten Kristiansen,
610 and Jun Wang. Snp detection for massively parallel whole-genome resequencing. *Genome*
611 *Research*, 19(6):1124–1132, 2009. doi:10.1101/gr.088013.108.
- 612 19 Li Fan, Pei Cao, J. Almeida, and A. Z. Broder. Summary cache: a scalable wide-area
613 web cache sharing protocol. *IEEE/ACM Transactions on Networking*, 8(3):281–293, 2000.
614 doi:10.1109/90.851975.
- 615 20 Leandro Lima, Blerina Sinimeri, Gustavo Sacomoto, H el ene Lopez-Maestre, Camille Marchet,
616 Vincent Miele, Marie-France Sagot, and Vincent Lacroix. Playing hide and seek with repeats
617 in local and global de novo transcriptome assembly of short RNA-seq reads. *Algorithms Mol*
618 *Biol*, 12, 2017.
- 619 21 Binghang Liu, Yujian Shi, Jianying Yuan, Xuesong Hu, Hao Zhang, Nan Li, Zhenyu Li,
620 Yanxiang Chen, Desheng Mu, and Wei Fan. Estimation of genomic characteristics by analyzing
621 k-mer frequency in de novo genome projects. *arXiv preprint arXiv:1308.2012*, 2013.
- 622 22 Yongchao Liu, Jan Schr oder, and Bertil Schmidt. Muskiet: a multistage k-mer spectrum-
623 based error corrector for Illumina sequence data. *Bioinformatics*, 29(3):308–315, 11 2012.
624 doi:10.1093/bioinformatics/bts690.
- 625 23 H el ene Lopez-Maestre, Lilia Brinza, Camille Marchet, Janice Kielbassa, Sylv ere Bastien,
626 Mathilde Boutigny, David Monnin, Adil El Filali, Claudia Marcia Carareto, Cristina Vieira,
627 Franck Picard, Natacha Kremer, Fabrice Vavre, Marie-France Sagot, and Vincent Lacroix.
628 SNP calling from RNA-seq data without a reference genome: identification, quantification,
629 differential analysis and impact on the protein sequence. *Nucleic Acids Research*, 44(19):e148–
630 e148, 2016. doi:10.1093/nar/gkw655.
- 631 24 Camille Marchet, Christina Boucher, Simon J Puglisi, Paul Medvedev, Mika el Salson, and
632 Rayan Chikhi. Data structures based on k-mers for querying large collections of sequencing
633 data sets. *Genome Research*, 31(1):1–12, 2021.
- 634 25 Guillaume Mar ais and Carl Kingsford. A fast, lock-free approach for efficient parallel
635 counting of occurrences of k-mers. *Bioinformatics*, 27(6):764–770, 01 2011. arXiv:https:
636 //academic.oup.com/bioinformatics/article-pdf/27/6/764/16902460/btr011.pdf, doi:
637 10.1093/bioinformatics/btr011.
- 638 26 Ali Mortazavi, Brian A Williams, Kenneth McCue, Lorian Schaeffer, and Barbara Wold.
639 Mapping and quantifying mammalian transcriptomes by rna-seq. *Nature Methods*, 5(7):621–
640 628, 2008. doi:10.1038/nmeth.1226.
- 641 27 Prashant Pandey, Michael A Bender, Rob Johnson, and Rob Patro. deBGR: an efficient and
642 near-exact representation of the weighted de Bruijn graph. *Bioinformatics*, 33(14):i133–
643 i141, 07 2017. arXiv:https://academic.oup.com/bioinformatics/article-pdf/33/14/
644 i133/25157228/btx261.pdf, doi:10.1093/bioinformatics/btx261.

- 645 28 Prashant Pandey, Michael A. Bender, Rob Johnson, and Rob Patro. A general-purpose
646 counting filter: Making every bit count. In *Proceedings of the 2017 ACM International*
647 *Conference on Management of Data*, SIGMOD '17, page 775–787, New York, NY, USA, 2017.
648 Association for Computing Machinery.
- 649 29 Prashant Pandey, Michael A Bender, Rob Johnson, and Rob Patro. Squeakr: an exact and
650 approximate k-mer counting system. *Bioinformatics*, 34(4):568–575, 10 2017. arXiv:<https://academic.oup.com/bioinformatics/article-pdf/34/4/568/25117209/btx636.pdf>, doi:
651 [10.1093/bioinformatics/btx636](https://doi.org/10.1093/bioinformatics/btx636).
652
- 653 30 Rob Patro, Geet Duggal, Michael I. Love, Rafael A. Irizarry, and Carl Kingsford. Salmon
654 provides fast and bias-aware quantification of transcript expression. *Nature methods*, 14(4):417–
655 419, 2017. URL: <https://pubmed.ncbi.nlm.nih.gov/28263959>.
- 656 31 Pavel A. Pevzner. 1-tuple dna sequencing: Computer analysis. *Journal of Biomolecular*
657 *Structure and Dynamics*, 7(1):63–73, 1989. PMID: 2684223. doi:[10.1080/07391102.1989.](https://doi.org/10.1080/07391102.1989.10507752)
658 [10507752](https://doi.org/10.1080/07391102.1989.10507752).
- 659 32 Rajeev Raman, Venkatesh Raman, and Srinivasa Rao Satti. Succinct indexable dictionaries
660 with applications to encoding k -ary trees, prefix sums and multisets. *ACM Trans. Algorithms*,
661 3(4):43, 2007. doi:[10.1145/1290672.1290680](https://doi.org/10.1145/1290672.1290680).
- 662 33 Yoshihiro Shibuya and Gregory Kucherov. Set-min sketch: a probabilistic map for power-law
663 distributions with application to k-mer annotation. *SeqBIM 2020*, 2020. doi:[10.1101/2020.](https://doi.org/10.1101/2020.11.14.382713)
664 [11.14.382713](https://doi.org/10.1101/2020.11.14.382713).
- 665 34 Daniel Dominic Sleator and Robert Endre Tarjan. A data structure for dynamic trees. *J.*
666 *Comput. Syst. Sci.*, 26(3):362–391, 1983. doi:[10.1016/0022-0000\(83\)90006-5](https://doi.org/10.1016/0022-0000(83)90006-5).
- 667 35 Raluca Uricaru, Guillaume Rizk, Vincent Lacroix, Elsa Quillery, Olivier Plantard, Rayan
668 Chikhi, Claire Lemaitre, and Pierre Peterlongo. Reference-free detection of isolated SNPs.
669 *Nucleic Acids Research*, 43(2):e11–e11, 11 2014. doi:[10.1093/nar/gku1187](https://doi.org/10.1093/nar/gku1187).
- 670 36 Reda Younsi and Dan MacLean. Using $2k + 2$ bubble searches to find single nucleotide
671 polymorphisms in k-mer graphs. *Bioinformatics*, 31(5):642–646, 10 2014. doi:[10.1093/](https://doi.org/10.1093/bioinformatics/btu706)
672 [bioinformatics/btu706](https://doi.org/10.1093/bioinformatics/btu706).
- 673 37 Birney E. Zerbino DR. Velvet: algorithms for de novo short read assembly using de Bruijn
674 graphs. *Genome Res.*, 18(5):821–9., 2008. PMID: 18349386. doi:[doi:doi:10.1101/gr.074492.107](https://doi.org/10.1101/gr.074492.107).