

# Techniques for Inverted Index Compression

GIULIO ERMANNIO PIBIRI, ISTI-CNR, Italy

ROSSANO VENTURINI, University of Pisa, Italy

The data structure at the core of large-scale search engines is the *inverted index*, which is essentially a collection of sorted integer sequences called *inverted lists*. Because of the many documents indexed by such engines and stringent performance requirements imposed by the heavy load of queries, the inverted index stores billions of integers that must be searched efficiently. In this scenario, *index compression* is essential because it leads to a better exploitation of the computer memory hierarchy for faster query processing and, at the same time, allows reducing the number of storage machines.

The aim of this article is twofold: first, surveying the encoding algorithms suitable for inverted index compression and, second, characterizing the performance of the inverted index through experimentation.

CCS Concepts: • **Information systems** → **Search index compression**; **Search engine indexing**.

Additional Key Words and Phrases: Inverted Indexes; Data Compression; Efficiency

## ACM Reference Format:

Giulio Ermanno Pibiri and Rossano Venturini. 2020. Techniques for Inverted Index Compression. *ACM Comput. Surv.* 1, 1 (August 2020), 35 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

## 1 INTRODUCTION

Consider a collection of textual documents each described, for this purpose, as a set of terms. For each distinct term  $t$  appearing in the collection, an integer sequence  $S_t$  is built and lists, in sorted order, all the identifiers of the documents (henceforth, docIDs) where the term appears. The sequence  $S_t$  is called the *inverted list*, or *posting list*, of the term  $t$  and the set of inverted lists for all the distinct terms is the subject of this article – the data structure known as the *inverted index*. Inverted indexes can store additional information about each term, such as the set of positions where the terms appear in the documents (in *positional indexes*) and the number of occurrences of the terms in the documents, i.e., their *frequencies* [14, 56, 113]. In this article we consider the docID-sorted version of the inverted index and we ignore additional information about each term. The inverted index is the data structure at the core of large-scale search engines, social networks and storage architectures [14, 56, 113]. In a typical use case, it is used to index millions of documents, resulting in several billions of integers. We mention some noticeable examples.

Classically, inverted indexes are used to support full-text search in databases [56]. Identifying a set of documents containing *all* the terms in a user query reduces to the problem of *intersecting* the inverted lists associated to the terms in the query. Likewise, an inverted list can be associated to a user in a social network (e.g., Facebook) and stores the sequence of all friend identifiers of the user [22]. Database systems based on SQL often precompute the list of row identifiers matching a specific frequent predicate over a large table, in order to speed up the execution of a query

---

Authors' addresses: Giulio Ermanno Pibiri, ISTI-CNR, Via Giuseppe Moruzzi 1, 56124, Pisa, Italy; Rossano Venturini, University of Pisa, Largo Bruno Pontecorvo 3, 56127, Pisa, Italy, [giulio.ermanno.pibiri@isti.cnr.it](mailto:giulio.ermanno.pibiri@isti.cnr.it), [rossano.venturini@unipi.it](mailto:rossano.venturini@unipi.it).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2020 Copyright held by the owner/author(s). Publication rights licensed to ACM.

0360-0300/2020/8-ART \$15.00

<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

Table 1. Timeline of techniques.

1949	Shannon-Fano [32, 93]	2005	Simple-9, Relative-10, and Carryover-12 [3]; RBUC [60]
1952	Huffman [43]	2006	PForDelta [114]; BASC [61]
1963	Arithmetic [1] <sup>1</sup>	2008	Simple-16 [112]; Tournament [100]
1966	Golomb [40]	2009	ANS [27]; Varint-GB [23]; Opt-PFor [111]
1971	Elias-Fano [30, 33]; Rice [87]	2010	Simple8b [4]; VSE [96]; SIMD-Gamma [91]
1972	Variable-Byte and Nibble [101]	2011	Varint-G8IU [97]; Parallel-PFor [5]
1975	Gamma and Delta [31]	2013	DAC [12]; Quasi-Succinct [107]
1978	Exponential Golomb [99]	2014	partitioned Elias-Fano [73]; QMX [103]; Roaring [15, 51, 53]
1985	Fibonacci-based [6, 37]	2015	BP32, SIMD-BP128, and SIMD-FastPFor [50]; Masked-VByte [84]
1986	Hierarchical bit-vectors [35]	2017	clustered Elias-Fano [80]
1988	Based on Front Coding [16]	2018	Stream-VByte [52]; ANS-based [63, 64]; Opt-VByte [83]; SIMD-Delta [104]; general-purpose compression libraries [77];
1996	Interpolative [65, 66]	2019	DINT [79]; Slicing [78]
1998	Frame-of-Reference (For) [39]; modified Rice [2]		
2003	SC-dense [11]		
2004	Zeta [8, 9]		

involving the conjunction of many predicates [42, 85]. Key-value storage is a popular database design principle, adopted by architectures such as Apache Ignite, Redis, InfinityDB, BerkeleyDB and many others. Common to all such architectures is the organization of data elements falling into the same bucket due to an hash collision: the list of all such elements is materialized, which is essentially an inverted list [24].

Because of the huge quantity of indexed documents and heavy query loads, *compressing* the inverted index is indispensable because it can introduce a twofold advantage over a non-compressed representation: feed faster memory levels with more data and, *hence*, speed up the query processing algorithms. As a result, the design of techniques that compress the index effectively while maintaining a noticeable decoding speed is a well-studied problem, that dates back to more than 50 years ago, and still a very active field of research. In fact, many representation for inverted lists are known, each exposing a different space/time trade-off: refer to the timeline shown in Table 1 and references therein.

**Organization.** We classify the techniques in a hierarchical manner by identifying *three* main classes. The first class consists of algorithms that compress a single integer (Section 2). The second class covers algorithms that compress many integers together, namely an inverted list (Section 3). The third class describes a family of algorithms that represent many lists together, i.e., the whole inverted index (Section 4). In our intention, this first part of the survey is devoted to readers who are new to the field of integer compression.

This hierarchical division is natural and intuitive. First, it reflects the flexibility of the algorithms, given that algorithms in a higher class can be used to represent the unit of compression of the algorithms in a lower class, but not the other way round. For example, an algorithm that compresses

<sup>1</sup>Actually somewhat before 1963. See Note 1 on page 61 in the book by Abramson [1].

a single integer at a time (first class) can obviously be used to represent a list of integers (unit of compression of the second class) by just concatenating the encoding of each single integer in the list. Second, it shows that less flexibility can be exploited to enlarge the “visibility” of the algorithms. For example, algorithms in the second class seek opportunities for better compression by looking for regions of similar integers in the list. Instead, algorithms in the third class seek for such regularities across many lists (or even across the whole index).

After the description of the techniques, we provide pointers to further readings (Section 5). The last part of the article is dedicated to the experimental comparison of the paradigms used to represent the inverted lists (Section 6). In our intention, this part is targeted to more experienced readers who are already familiar with the research field and the practical implementation of the techniques. We release the full experimental suite at [https://github.com/jermp/2i\\_bench](https://github.com/jermp/2i_bench), in the hope of spurring further research in the field. We conclude the survey by summarizing experimental lessons and discussing some future research directions (Section 7).

## 2 INTEGER CODES

The algorithms we consider in this section compress a single integer. We first introduce some preliminary notions that help to better illustrate such algorithms. Let  $x > 0$  indicate the integer to be represented<sup>2</sup>. Sometimes we assume that an upper bound on the largest value that  $x$  can take is known, the so-called “universe” of representation, and we indicate that with  $U$ . Therefore, it holds  $x \leq U, \forall x$ . Let  $C(x)$  be the bit string encoding  $x$  – the *codeword* of  $x$  according to the code  $C$  – and  $|C(x)|$  its length in bits.

The most classical solution to the problem of integer coding is to assign  $x$  a *uniquely-decodable* variable-length code, in order to decode without ambiguity (thus, correctly) from left to right. The aim of the algorithms we present in the following is to assign the smallest codeword as possible. In this regard, a distinction is made between the identification of a *set of codeword lengths* and the subsequent *codeword assignment* phase, that is the mapping from integer identifiers to binary strings. Once the codeword lengths have been computed, the specific codewords are actually irrelevant provided that no codeword is a prefix of another one (prefix-free condition) in order to guarantee unique decodability – a key observation widely documented in the literature [49, 59, 67, 68, 109]. Therefore, we can think of a code as being a (multi-) set of codeword lengths. This interpretation has the advantage of being independent from the specific codewords and makes possible to choose the assignment that is best suited for fast decoding – an important property as far as practicality is concerned for inverted index compression. Throughout this section we opt for a *lexicographic assignment* of the codewords, that is the codewords are in the same lexicographic order as the integers they represent. This property will be exploited for fast decodability in Section 2.1. The crucial fact about these initial remarks is that *all* prefix-free codes can be arranged in this way, hence providing a “canonical” interface for their treatment.

Another key notion is represented by the *Kraft–McMillan inequality* [47, 57] that gives a necessary and sufficient condition for the existence of a uniquely-decodable code for a given set of codeword lengths, where no codeword is a prefix of another one. More formally, it must hold  $\sum_{x=1}^u 2^{-|C(x)|} \leq 1$  for the code to be uniquely-decodable.

It should be intuitive that no code is optimal for *all* possible integer distributions. According to Shannon [93], the ideal codeword length of the integer  $x$  should be  $\log_2(1/\mathbb{P}(x))$  bits long, where  $\mathbb{P}(x)$  is the probability of occurrence of  $x$  in the input. Therefore, by solving the equation

<sup>2</sup>Throughout the section we present the codes for positive integers. Also, we use 1-based indexes for arrays. This is a deliberate choice for illustrative purposes. The reader should be aware that some implementations at [https://github.com/jermp/2i\\_bench](https://github.com/jermp/2i_bench) may differ from the content of the paper and should be careful in comparing them.

$|C(x)| = \log_2(1/\mathbb{P}(x))$  with respect to  $\mathbb{P}(x)$ , the distribution for which the considered integer code is optimal can be derived.

Lastly, we also remark that sometimes it could be useful to implement a sub-optimal code if it allows faster decoding for a given application, and/or to organize the output bit stream in a way that is more suitable for special hardware instructions such as SIMD (Single-Instruction-Multiple-Data) [20]. SIMD is a computer organization that exploits the independence of multiple data objects to execute a single instruction on these objects simultaneously. Specifically, a single instruction is executed for every element of a *vector* – a large(r) machine register that packs multiple elements together, for a total of 128, 256, or even 512 bits of data. SIMD is widely used to accelerate the execution of many data-intensive tasks, and (usually) an optimizing compiler is able to automatically “vectorize” code snippets to make them run faster. Many algorithms that we describe in this article exploit SIMD instructions.

## 2.1 Encoding and decoding prefix-free codes

We now illustrate how the lexicographic ordering of the codewords can be exploited to achieve efficient encoding and decoding of prefix-free codes. By “efficiency” we mean that a small and fixed number of instructions is executed for each encoded/decoded integer, thus avoiding the potentially expensive bit-by-bit processing. The content of this section is based on the work by Moffat and Turpin [67] that also laid the foundations for (part of) Chapter 4 of their own book [68] and the descriptions in Section 3.1 and 3.2 of the survey by Moffat [59].

We are going to refer to Table 2a as an example code. The table shows the codewords assigned to the integers 1..8 by the *gamma* code, the first non-trivial code we will describe later in Section 2.3. Let  $M$  be the length of the longest codeword, that is  $M = \max_x \{|C(x)|\}$ . In our example,  $M = 7$ . The column headed *lengths* reports the lengths (in bits) of the codewords, whereas *values* are the decimal representation of the codewords seen as left-justified  $M$ -bit integers. If such columns are represented with two parallel arrays indexed by the  $x$  value, then the procedures for encoding/decoding are easily derived as follows. To encode  $x$ , just write to the output stream the *most significant* (from the left)  $lengths[x]$  bits of  $values[x]$ . To decode  $x$ , we use a variable *buffer* always holding  $M$  bits from the input stream. The value assumed by *buffer* is, therefore, an  $M$ -bit integer that we search in the array *values* determining the codeword length  $\ell = length[x]$  such that  $values[x] \leq buffer < values[x+1]$ . Now that  $\ell$  bits are consumed from *buffer*, other  $\ell$  bits are fetched from the input stream via a few masking and shifting operations.

For example, to encode the integer 4, we omit the  $lengths[4] = 5$  most significant bits from the binary representation of  $values[4] = 96$  as a 7-bit integer, that is 1100000. The obtained codeword for 4 is, therefore, 11000. Instead, assume we want to decode the next integer from a *buffer* configuration of 1010100. This 7-bit integer is 84. By searching 84 in the *values* array, we determine the index  $x = 3$  as  $80 = values[3] < 84 < values[4] = 96$ . Therefore, the decoded integer is 3 and we can fetch the next  $lengths[3] = 3$  bits from the input stream. (It is easy to see that the *buffer* configuration 1010100 holds the encoded values 3, 0, and 2.)

However, the cost of storing the *lengths* and *values* arrays can be large because they can hold as many values as  $U$ . The universe  $U$  is typically in the range of tens of millions for a typical inverted index benchmark (see also Table 10 at page 10 for a concrete example), thus resulting in large lookup tables that do not fit well in the computer cache hierarchy. Fortunately enough, it is possible to replace such arrays with two compact arrays of just  $M + 1$  integers each when the codewords are assigned in lexicographic order. Recall that we have defined  $M$  to be the longest codeword length. (In our example from Table 2,  $M$  is 7.) We expect to have  $M \leq U$  and, in particular,  $M \ll U$ , which is always valid in practice unless  $U$  is very small.

Table 2. In (a), an example prefix-free code for the integers 1..8, along with associated codewords, codeword lengths and corresponding left-justified, 7-bit, integers. The codewords are left-justified to better highlight their lexicographic order. In (b), the compact version of the table in (a), used by the encoding/decoding procedures coded in Fig. 1. The “values” and “first” columns are padded with a sentinel value (in gray) to let the search be well defined.

(a)				(b)		
$x$	codewords	lengths	values	lengths	first	values
1	0	1	0	1	1	0
2	100	3	64	2	2	64
3	101	3	80	3	2	64
4	11000	5	96	4	4	96
5	11001	5	100	5	4	96
6	11010	5	104	6	8	112
7	11011	5	108	7	8	112
8	1110000	7	112	–	9	127
–	–	–	127			

```

1 Encode( $x$ ) :
2   determine  $\ell$  such that  $first[\ell] \leq x < first[\ell + 1]$ 
3    $offset = x - first[\ell]$ 
4    $jump = 1 \ll (M - \ell)$ 
5   Write( $(values[\ell] + offset \times jump) \gg (M - \ell), \ell$ )

1 Decode() :
2   determine  $\ell$  such that  $values[\ell] \leq buffer < values[\ell + 1]$ 
3    $offset = (buffer - values[\ell]) \gg (M - \ell)$ 
4    $buffer = ((buffer \ll \ell) \& MASK) + Take(\ell)$ 
5   return  $first[\ell] + offset$ 

```

► MASK is the constant  $2^M - 1$ .

Fig. 1. Encoding and decoding procedures using two parallel arrays *first* and *values* of  $M + 1$  values each.

Table 2b shows the “compact” version of Table 2a, where other two arrays, *first* and *values*, of  $M + 1$  integers each are used. Both arrays are now indexed by codeword length  $\ell$ . In particular,  $first[\ell]$  is the first integer that is assigned a codeword of length equal to  $\ell$ , with  $values[\ell]$  being the corresponding  $M$ -bit integer representation. For example,  $first[3] = 2$  because 2 is the first integer represented with a codeword of length 3. Note that not every possible codeword length could be used. In our example, we are not using codewords of length 2, 4 and 6. These unused lengths generate some “holes” in the *first* and *values* arrays. A hole at position  $\ell$  is then filled with the value corresponding to the smallest codeword length  $\ell' > \ell$  that is used by the code. For example, we have a hole at position  $\ell = 2$  because we have no codeword of length 2. Therefore,  $first[2]$  and  $values[2]$  are filled with 2 and 64 respectively because these are the values corresponding to the codeword length  $\ell' = 3$ . With these two arrays it is possible to derive the compact pseudo code illustrated in Fig. 1, whose correctness is left to the reader. The function *Write*(*val*, *len*) writes the *len* low bits of the value *val* to the output stream; conversely, the function *Take*(*len*) fetches the next *len* bits from the input stream and interprets them as an integer value. We now discuss some

salient aspects of the pseudo code along with examples, highlighting the benefit of working with lexicographically-sorted codewords.

Assigning lexicographic codewords makes possible to use the offsets computed in line 3 of both listings in Fig. 1 to perform encoding/decoding, essentially allowing a sparse representation of the mapping from integers to codewords and vice versa. As an example, consider the encoding of the integer 6. By searching the *first* array, we determine  $\ell = 5$ . Now, the difference  $offset = 6 - first[5] = 6 - 4 = 2$  indicates that 6 is the  $(offset + 1)$ -th integer, the 3rd in this case, that is assigned a codeword of length 5. Therefore, starting from  $values[5] = 96$  we can derive the  $M$ -bit integer corresponding to the encoding of 6 using such offset, via  $96 + 2 \times 4 = 104$ . It is easy to see that this computation is correct only when the codewords are assigned lexicographically, otherwise it would not be possible to derive the  $M$ -bit integer holding the codeword. The same reasoning applies when considering the decoding algorithm.

Another advantage of working with ordered codewords is that both *first* and *values* arrays are sorted, thus binary search can be employed to implement the identification of  $\ell$  at line 2 of both encoding and decoding algorithms. This step is indeed the most expensive one. Linear search could result faster than binary search for its cache-friendly nature, especially when  $M$  is small. If we are especially concerned with decoding efficiency and can trade a bit more working space, it is possible to identify  $\ell$  via direct addressing, i.e., in  $O(1)$  per decoded symbol, using a  $2^M$ -element table indexed by *buffer*. Other options are possible, including hybrid strategies using a blend of search and direct addressing: for a discussion of these options, again refer to the work by Moffat and Turpin [67, 68], Moffat [59, Section 3.1 and 3.2], and references therein.

## 2.2 Unary and Binary

Perhaps the most primitive form of integer representation is *unary* coding, that is the integer  $x$  is represented by a run of  $x - 1$  ones plus a final zero:  $1^{x-1}0$ . The presence of the final 0 bit implies that the encoding is uniquely-decodable: decoding bit-by-bit will keep reading ones until we hit a 0 and then report the number of read ones by summing one to this quantity. Because we need  $x$  bits to represent the integer  $x$ , that is  $|U(x)| = x$ , this encoding strongly favours small integers. For example, we can represent 2 with just 2 bits (10) but we would need 503 bits to represent the integer 503. Solving  $\lceil \log_2(1/\mathbb{P}(x)) \rceil = x$  yields that the unary code is optimal whenever  $\mathbb{P}(x) = 1/2^x$ .

We indicate with  $\text{bin}(x, k)$  the *binary representation* of an integer  $0 \leq x < 2^k$  using  $k$  bits. When we just write  $\text{bin}(x)$ , it is assumed that  $k = \lceil \log_2(x + 1) \rceil$ , which is the minimum number of bits necessary to represent  $x$ . We also distinguish between  $\text{bin}(x, k)$  and the *binary codeword*  $B(x)$  assigned to an integer  $x > 0$ . Given that we consider positive integers only, we use the convention that  $B(x)$  is  $\text{bin}(x - 1)$  throughout this section. For example,  $B(6) = \text{bin}(5) = 101$ . See the second column of Table 3 for more examples. This means that, for example,  $|B(4)|$  is 2 and not 3;  $B(503)$  just needs  $\lceil \log_2 503 \rceil = 9$  bits instead of 503 as needed by its unary representation. The problem of binary coding is that it is *not* uniquely-decodable unless we know the number of bits that we dedicate to the representation of *each* integer in the coded stream. For example, if the integers in the stream are drawn from a universe  $U$  bounded by  $2^k$  for some  $k > 0$ , then each integer can be represented with  $\lceil \log_2 U \rceil \leq k$  bits, with an implied distribution of  $\mathbb{P}(x) = 1/2^k$ . (Many compressors that we present in Section 3 exploit this simple strategy.) If  $U = 2^k$ , then the distribution simplifies to  $\mathbb{P}(x) = 1/U$  (i.e., *uniform*).

The following definition will be useful. Consider  $x \in [0, b - 1]$  for some  $b > 0$  and let  $c = \lceil \log_2 b \rceil$ . We define the *minimal* binary codeword assigned to  $x$  in the interval  $[0, b - 1]$  as  $\text{bin}(x, c - 1)$  if  $x < 2^c - b$ ,  $\text{bin}(x + 2^c - b, c)$  otherwise. Note that if  $b$  is a power of two the minimal binary codeword for  $x$  is  $\text{bin}(x, c)$ .



Table 3. The integers 1..8 as represented with several codes. The “.” symbol highlights the distinction between different parts of the codes and has a purely illustrative purpose: it is not included in the final coded representation.

$x$	$U(x)$	$B(x)$	$\gamma(x)$	$\delta(x)$	$G_2(x)$	$\text{Exp}G_2(x)$	$Z_2(x)$
1	0	0	0.	0.	0.0	0.00	0.0
2	10	1	10.0	100.0	0.1	0.01	0.10
3	110	10	10.1	100.1	10.0	0.10	0.11
4	1110	11	110.00	101.00	10.1	0.11	10.000
5	11110	100	110.01	101.01	110.0	10.000	10.001
6	111110	101	110.10	101.10	110.1	10.001	10.010
7	1111110	110	110.11	101.11	1110.0	10.010	10.011
8	11111110	111	1110.000	11000.000	1110.1	10.011	10.1000

### 2.3 Gamma and Delta

The two codes we now describe were introduced by Elias [31] and are called *universal* because the length of these codes is  $O(\log x)$  bits for every integer  $x$ , thus a constant factor away from the optimal binary representation of length  $|\text{bin}(x)| = \lceil \log_2(x+1) \rceil$  bits. Additionally, they are uniquely-decodable.

The *gamma* code for  $x$  is made by the unary representation of  $|\text{bin}(x)|$  followed by the  $|\text{bin}(x)|-1$  least significant bits of  $\text{bin}(x)$ . Therefore,  $|\gamma(x)| = 2|\text{bin}(x)| - 1$  bits and  $\mathbb{P}(x) \approx 1/(2x^2)$ . Bit-by-bit decoding of  $\gamma(x)$  is simple too. First read the unary code, say  $\ell$ . Then sum to  $2^{\ell-1}$  the integer represented by the next  $\ell - 1$  bits. For example, the integer 113 is represented as  $\gamma(113) = 1111110.110001$ , because  $\text{bin}(113) = 1110001$  is 7 bits long.

The key inefficiency of the gamma code lies in the use of the unary code for the representation of  $|\text{bin}(x)|$ , which may become very large for big integers. To overcome this limitation, the *delta* code replaces  $U(|\text{bin}(x)|)$  with  $\gamma(|\text{bin}(x)|)$  in the  $\delta$  representation of  $x$ . The number of bits required by  $\delta(x)$  is, therefore,  $|\gamma(|\text{bin}(x)|)| + |\text{bin}(x)| - 1$ . The corresponding distribution is  $\mathbb{P}(x) \approx 1/(2x(\log_2 x)^2)$ . Bit-by-bit decoding of  $\delta$  codes follows automatically from that of  $\gamma$  codes. Again, the integer 113 is represented as  $\delta(113) = 11011.110001$ . The first part of the encoding, 11011, is the  $\gamma$  representation of 7, which is the length of  $\text{bin}(113)$ . Table 3 shows the integers 1..8 as encoded with  $\gamma$  and  $\delta$  codes.

In order to decode gamma codes faster on modern processors, a simple variant of gamma is proposed by Schlegel et al. [91] and called *k-gamma*. Groups of  $k$  integers are encoded together, with  $k = 2, 3$  or  $4$ , using the same number of bits. Thus, instead of recording the unary length for each integer, only the length of the largest integer in the group is written. This leads to a higher compression ratio if the  $k$  integers are close to each other, namely they require the same codeword length. On the other hand, if the largest integers in the group requires more bits than the other integers, this encoding is wasteful compared to the traditional gamma. However, decoding is faster: once the binary length has been read, a group of  $k$  integers is decoded in parallel using SIMD instructions. Similarly, Trotman and Lilly [104] introduced a SIMD version of delta codes. A 512-bit payload is broken down into its  $16 \times 32$ -bit integers and the base-2 magnitude of the largest integer is written using gamma coding as a selector (writing the selector in unary code gives 16-gamma). Although not as fast as *k-gamma*, the representation is faster to decode compared to decoding bit by bit.

## 2.4 Golomb

In 1966 Golomb introduced a parametric code that is a hybrid between unary and binary [40]. The Golomb code of  $x$  with parameter  $b > 1$ ,  $G_b(x)$ , consists in the representation of two pieces: the *quotient*  $q = \lfloor (x - 1)/b \rfloor$  and the *remainder*  $r = x - q \times b - 1$ . The quantity  $q + 1$  is written in unary, followed by a minimal binary codeword assigned to  $r \in [0, b - 1]$ . Clearly, the closer  $b$  is to the value of  $x$  the smaller the value of  $q$ , with consequent better compression and faster decoding speed. Table 3 shows an example code with  $b = 2$ . Let us consider the code with  $b = 5$ , instead. From the definition of minimal binary codeword, we have that  $c = \lceil \log_2 5 \rceil = 3$  and  $2^c - b = 3$ . Thus the first 3 reminders, 0..2, are always assigned the first 2-bit codewords 00, 01 and 10 respectively. The reminders 3 and 4 are instead assigned codewords 110 and 111 as given by  $\text{bin}(3 + 3, 3)$  and  $\text{bin}(4 + 3, 3)$ , respectively. Decoding just reverts the encoding procedure. After the unary prefix, always read  $c - 1$  bits, with  $c = \lceil \log_2 b \rceil$ . If these  $c - 1$  bits give a quantity  $r$  that is less than  $2^c - b$ , then stop and work with the reminder  $r$ . Instead, if  $r \geq 2^c - b$  then fetch another bit and compute  $r$  as the difference between this  $c$ -bit number and the quantity  $2^c - b$ .

Golomb was the first to observe that if  $n$  integers are drawn *at random* from a universe of size  $U$ , then the gaps between the integers follow a geometric distribution  $\mathbb{P}(x) = p(1 - p)^{x-1}$  with parameter  $p = n/U$  being the probability to find an integer  $x$  among the ones selected. It is now clear that the optimal value for  $b$  depends on  $p$  and it can be shown that this value is the integer closest to  $-1/\log_2(1 - p)$ , i.e., the value that satisfies  $(1 - p)^b \approx 1/2$ . Doing the math we have  $b \approx 0.69/p$ , which is a good approximation of the optimal value and can be used to define a Golomb code with parameter  $b$ . This code is optimal for the geometric distribution  $\mathbb{P}(x) = p(1 - p)^{x-1}$ . Gallager and Van Voorhis [38] showed that the optimal value for  $b$  can be computed as  $b = -\lceil \log_2(2 - p)/\log_2(1 - p) \rceil$ .

## 2.5 Rice

The Rice code [86, 87] is actually a special case of the Golomb code for which  $b$  is set to  $2^k$ , for some  $k > 0$  (sometimes also referred to as the Golomb-Rice code). Let  $\text{Rice}_k(x)$  be the Rice code of  $x$  with parameter  $k > 0$ . In this case the remainder  $r$  is always written in  $k$  bits. Therefore, the length of the Rice code is  $|\text{Rice}_k(x)| = \lfloor (x - 1)/2^k \rfloor + k + 1$  bits. To compute the optimal parameter for the Rice code, we just pretend to be constructing an optimal Golomb code with parameter  $b$  and then find two integers,  $l$  and  $r$ , such that  $2^l \leq b \leq 2^r$ . One of these two integers will be the optimal value of  $k$  for the Rice code. (We also point the interested reader to the technical report by Kiely [46] for a deep analysis about the optimal values of the Rice parameter.)

## 2.6 Exponential Golomb

The *exponential* Golomb code proposed by Teuhola [99] logically defines a vector of “buckets”

$$B = \left[ 0, 2^k, \sum_{i=0}^1 2^{k+i}, \sum_{i=0}^2 2^{k+i}, \sum_{i=0}^3 2^{k+i}, \dots \right], \text{ for some } k \geq 0$$

and encodes an integer  $x$  as a bucket identifier plus an offset relative to the bucket. More specifically, the code  $\text{ExpG}_k(x)$  is obtained as follows. We first determine the bucket where  $x$  belongs to, i.e., the index  $h \geq 1$  such that  $B[h] < x \leq B[h + 1]$ . Then  $h$  is coded in unary, followed by a minimal binary codeword assigned to  $x - B[h] - 1$  in the shrunk interval  $[0, B[h + 1] - B[h] - 1]$ . (Since  $\log_2(B[h + 1] - B[h])$  is always a power of 2 for the choice of  $B$  above, the binary codeword of  $x$  is  $\text{bin}(x - B[h] - 1, \log_2(B[h + 1] - B[h]))$ ).

Table 3 shows an example for  $k = 2$ . Note that  $\text{ExpG}_0$  coincides with Elias'  $\gamma$ .



Table 4. The integers 1..8 as represented with Fibonacci-based codes. In (a), the final control bit is highlighted in bold font and the relevant Fibonacci numbers  $F_i$  involved in the representation are also shown at the bottom of the table. In (b), the “canonical” lexicographic codewords are presented.

(a) “original” codewords

$x$	$F(x)$					
1	1	<b>1</b>				
2	0	1	<b>1</b>			
3	0	0	1	<b>1</b>		
4	1	0	1	<b>1</b>		
5	0	0	0	1	<b>1</b>	
6	1	0	0	1	<b>1</b>	
7	0	1	0	1	<b>1</b>	
8	0	0	0	0	1	<b>1</b>
$F_i$	1	2	3	5	8	13

(b) lexicographic codewords

$x$	$F(x)$						
1	0	0					
2	0	1	0				
3	0	1	1	0			
4	0	1	1	1			
5	1	0	0	0	0		
6	1	0	0	0	1		
7	1	0	0	1	0		
8	1	0	0	1	1	0	

## 2.7 Zeta

Boldi and Vigna [8, 9] introduced the family of *zeta* codes that is optimal for integers distributed according to a power law with small exponent  $\alpha$  (e.g., less than 1.6), that is  $\mathbb{P}(x) = 1/(\zeta(\alpha)x^\alpha)$ , where  $\zeta(\cdot)$  denotes the Riemann zeta function. The zeta code  $Z_k$  is an exponential Golomb code relative to a vector of “buckets”  $[0, 2^k - 1, 2^{2k} - 1, 2^{3k} - 1, \dots]$ . Again, Table 3 shows an examples for  $k = 2$ . Note that  $Z_1$  coincides with  $\text{ExpG}_0$ , therefore also  $Z_1$  is identical to Elias’  $\gamma$ .

For example, let us consider  $Z_2(5)$ . The value of  $h$  is 2 because  $2^2 - 1 < 5 < 2^4 - 1$ . Therefore the first part of the code is the unary representation of 2. Now we have to assign a minimal binary codeword to  $5 - (2^2 - 1) - 1 = 1$  using 3 bits, that is  $\text{bin}(1, 3) = 001$ . A more involved example is the one for, say,  $Z_3(147)$ . In this case, we have  $h = 3$ , thus the interval of interest is  $[2^6 - 1, 2^9 - 1]$ . Now we have to assign  $147 - (2^6 - 1) - 1 = 83$  a minimal binary codeword in the interval  $[0, 448]$ . Since 83 is more than the left extreme  $2^6 - 1$ , we have to write  $\text{bin}(147, 9)$  for a final codeword of  $110.010010011$ .

## 2.8 Fibonacci

Fraenkel and Klein [37] introduced in 1985 a class of codes based on *Fibonacci* numbers [71] and later generalized by Apostolico and Fraenkel [6]. The encoding is a direct consequence of the *Zeckendorf’s theorem*: every positive integer can be uniquely represented as the sum of some, non adjacent, Fibonacci numbers. Let  $F_i = F_{i-1} + F_{i-2}$  define the  $i$ -th Fibonacci number for  $i > 2$ , with  $F_1 = 1$  and  $F_2 = 2$ . Then we have:  $F_3 = 3, F_4 = 5, F_5 = 8, F_6 = 13$ , etc. The Fibonacci encoding  $F(x)$  of an integer  $x$  is obtained by: (1) emitting a 1 bit if the  $i$ -th Fibonacci number is used in the sum giving  $x$ , or emitting a 0 bit otherwise; (2) appending a final control 1 bit to ensure unique decodability. Table 4a shows the first 8 integers as encoded with this procedure, where we highlighted in bold font the final control bit. For example,  $7 = F_2 + F_4$ , thus  $F(7)$  will be given by 4 bits where the second and the fourth are 1, i.e., 0101, plus the control 1 bit for a final codeword of 01011.

Note that the codewords assigned by the procedure described above are not lexicographically-sorted in the integers they represent. However, if we first compute the codeword lengths we can then generate a set of lexicographically-sorted codewords in a rather simple way, therefore obtaining a Fibonacci-based code that can be encoded/decoded with the procedures we have illustrated in

Section 2.1. Given a non-decreasing sequence of codeword lengths  $[\ell_1, \dots, \ell_n]$  satisfying the Kraft-McMillan inequality (see the beginning of Section 2), the corresponding codewords are generated as follows. The first codeword is always the bit string of length  $\ell_1$  that is  $0^{\ell_1}$ . Now, let  $\ell = \ell_1$ . For all  $i = 2, \dots, n$  we repeat the following two steps.

- (1) Let  $c$  be the next lexicographic codeword of  $\ell$  bits. If  $\ell_i = \ell$ , then we just emit  $c$ . Otherwise,  $c$  is padded with possible  $0$  bits to the right until we have a  $\ell_i$ -bit codeword.
- (2) We set  $\ell = \ell_i$ .

Note that the way we define  $c$  in step (1) guarantees that the generated code is prefix-free.

For our example in Table 4b, the sequence of codeword lengths is  $[2, 3, 4, 4, 5, 5, 6]$ . Let us generate the first 4 codewords. The first codeword is therefore  $00$ , with  $\ell_1 = 2$ . The next codeword length  $\ell_2$  is 3, thus we pad the next 2-bit codeword following  $00$ , i.e.,  $01$ , with a  $0$  and obtain the 3-bit codeword  $010$ . The next codeword length is 4, thus we obtain the codeword  $0110$ . The next codeword length is 4 again and the codeword is just obtained by assigning the codeword following  $0110$  in lexicographic order, that is  $0111$ .

There is a closed-form formula for computing the  $i$ -th Fibonacci number,  $i \geq 1$ , called Binet's formula:

$$F_i = \frac{1}{\sqrt{5}} \left[ \left( \frac{1 + \sqrt{5}}{2} \right)^{i+1} - \left( \frac{1 - \sqrt{5}}{2} \right)^{i+1} \right] \approx \left( \frac{1 + \sqrt{5}}{2} \right)^{i+1} = \phi^{i+1},$$

where  $\phi = \frac{1+\sqrt{5}}{2}$  is the so-called *golden ratio*. Using this formula, it can be shown that the codeword length of  $F(x)$  is approximately equal to  $1 + \log_\phi x$  bits. Therefore, the corresponding distribution is  $\mathbb{P}(x) = 1/(2x^{1/\log_2 \phi}) \approx 1/(2x^{1.44})$ . This implies that Fibonacci-based codes are shorter than  $\gamma$  for  $x > 3$ ; and as good as or even better than  $\delta$  for a wide range of practical values ranging from  $F_2 = 2$  to  $F_{19} = 6765$ .

## 2.9 Variable-Byte

The codes described in the previous sections are *bit-aligned* as they do not represent an integer using a multiple of a fixed number of bits, e.g., a byte. But reading a stream of bits in chunks where each chunk is a byte of memory (or a multiple of a byte, e.g., a memory word – 4 or 8 bytes), is simpler and faster because the data itself is written in memory in this way. Therefore, it could be preferable to use *byte-aligned* or *word-aligned* codes when decoding speed is the main concern rather than compression effectiveness.

*Variable-Byte*, first described by Thiel and Heaps [101], is the most popular and simplest byte-aligned code: the binary representation of a non-negative integer is split into groups of 7 bits which are represented as a sequence of bytes. In particular, the 7 least significant bits of each byte are reserved for the data whereas the most significant, called the *continuation bit*, is equal to 1 to signal continuation of the byte sequence. The last byte of the sequence has its 8-th bit set to 0 to signal, instead, the termination of the byte sequence. The main advantage of Variable-Byte codes is decoding speed: we just need to read one byte at a time until we found a value smaller than  $2^7$ . Conversely, the number of bits to encode an integer cannot be less than 8, thus Variable-Byte is only suitable for large numbers and its compression ratio may not be competitive with the one of bit-aligned codes for small integers. Variable-Byte uses  $\lceil \lceil \log_2(x+1) \rceil / 7 \rceil \times 8$  bits to represent the integer  $x$ , thus it is optimal for the distribution  $\mathbb{P}(x) \approx \sqrt[7]{1/x^8}$ . For example, the integer 65,790 is represented as **000001**00.1**0000001**.11111110, where we mark the control bits in bold font. Also notice the padding bits in the first byte starting from the left, inserted to align the binary representation of the number to a multiple of 8 bits.

Table 5. The integers 1..20 as represented by SC(4, 4)- and SC(5, 3)-dense codes respectively.

$x$	SC(4, 4, $x$ )	SC(5, 3, $x$ )	$x$	SC(4, 4, $x$ )	SC(5, 3, $x$ )
1	000	000	11	101.010	110.000
2	001	001	12	101.011	110.001
3	010	010	13	110.000	110.010
4	011	011	14	110.001	110.011
5	100.000	100	15	110.010	110.100
6	100.001	101.000	16	110.011	111.000
7	100.010	101.001	17	111.000	111.001
8	100.011	101.010	18	111.001	111.010
9	101.000	101.011	19	111.010	111.011
10	101.001	101.100	20	111.011	111.100

*Nibble* coding is a simple variation of this strategy where 3 bits are used for data instead of 7, which is optimal for the distribution  $\mathbb{P}(x) \approx \sqrt[3]{1/x^4}$ .

Culpepper and Moffat [21] describe a byte-aligned code with the property that the first byte of each codeword defines the length of the codeword, which makes decoding simpler and faster.

Various enhancements were proposed to accelerate the sequential decoding speed of Variable-Byte. For example, in order to reduce the probability of a branch misprediction that leads to higher throughput and helps keeping the CPU pipeline fed with useful instructions, the control bits can be grouped together. If we assume that the largest represented integer fits into four bytes, we have to distinguish between only four different byte-lengths, thus two bits are sufficient. In this way, groups of four integers require one control byte only. This optimization was introduced in Google's Varint-GB format [23], which is faster to decode than the original Variable-Byte code.

Working with byte-aligned codes also opens the possibility of exploiting the parallelism of SIMD instructions to further enhance the sequential decoding speed. This is the approach taken by the proposals Varint-G8IU [97], Masked-VByte [84] and Stream-VByte [52] that we overview below.

Varint-G8IU [97] uses a format similar to the one of Varint-GB: one control byte describes a variable number of integers in a data segment of exactly eight bytes, therefore each group can contain between two and eight compressed integers. Masked-VByte [84] works, instead, directly on the original Variable-Byte format. The decoder first gathers the most significant bits of consecutive bytes using a dedicated SIMD instruction. Then using previously-built look-up tables and a shuffle instruction, the data bytes are permuted to obtain the decoded integers. Stream-VByte [52] separates the encoding of the control bits from the data bits by writing them into separate streams. This organization permits to decode multiple control bits simultaneously and, consequently, to reduce data dependencies that can stop the CPU pipeline execution when decoding the data stream.

## 2.10 SC-Dense

In Variable-Byte encoding the value  $2^7$  acts as a separator between *stoppers*, i.e., all values in  $[0, 2^7 - 1]$ , and *continuers*, i.e., all values in  $[2^7, 2^8 - 1]$ . A generalization of the encoding can be obtained by changing the separator value, thus enlarging or restricting the cardinalities of the set of continuers and stoppers. In general, the values from 0 to  $c - 1$  are reserved to stoppers and the values from  $c$  to  $c + s - 1$  to continuers, provided that  $c + s = 2^8$ . Intuitively, changing the separating value can better adapt to the distribution of the integers to be encoded. For example, if most integers are larger than (say) 127, then it is convenient to have more continuers. This is the main idea behind the SC-dense code introduced by Brisaboa et al. [11].

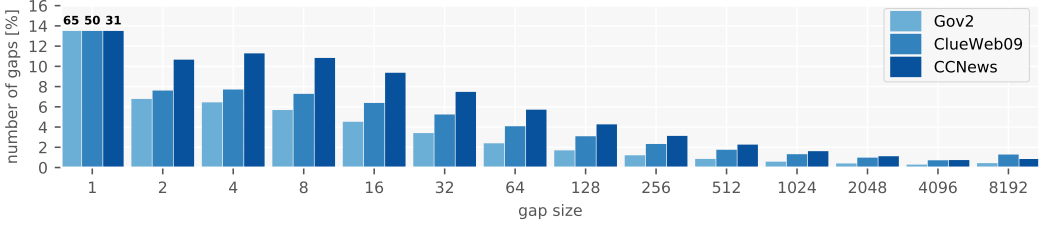


Fig. 2. Distribution of the gaps in the real-world datasets Gov2, ClueWeb09, and CCNews.

Given the integer  $x$ , its  $SC(s, c, x)$  representation is obtained as follows. Let  $k(x) \geq 1$  be the number of  $\lceil \log_2(s + c) \rceil$ -bit words needed by the representation of  $x$ . This value  $k(x)$  will be such that

$$s \frac{c^{k(x)-1} - 1}{c - 1} \leq x < s \frac{c^{k(x)} - 1}{c - 1}.$$

If  $k(x) = 1$ , the representation is just the stopper  $x - 1$ . Otherwise, let  $y = \lfloor (x - 1)/s \rfloor$  and  $x' = x - (sc^{k(x)-1} - s)/(c - 1)$ . In this case,  $x > s$  and  $k(x)$  is given by  $\lfloor (y - 1)/c \rfloor + 2$ . The representation is given by a sequence of  $k(x) - 1$  continers, made by repeating the continuer  $s + c - 1$  for  $k(x) - 2$  times followed by the continuer  $s + ((y - 1) \bmod c)$ , plus the final stopper  $(x' - 1) \bmod s$ . The number of bits required by encoding of  $x$  is  $k(x) \lceil \log_2(s + c) \rceil$ , thus it follows that  $\mathbb{P}(x) \approx (s + c)^{-k(x)}$ . It is also possible to compute via dynamic programming the optimal values for  $s$  and  $c$  given the probability distribution of the integers [11].

Table 5 shows the codewords assigned to the integers 1..20 by the dense codes  $SC(4, 4)$  and  $SC(5, 3)$  respectively. Let us consider the encoding of  $x = 13$  under the code  $SC(5, 3)$ . In this example, we have  $y = \lfloor (13 - 1)/5 \rfloor = 2$  and  $k(13) = \lfloor (2 - 1)/3 \rfloor + 2 = 2$ . The only continuer is therefore given by  $5 + ((2 - 1) \bmod 3) = 6$ , i.e., 110. Since  $x' = 3$ , the stopper is  $(3 - 1) \bmod 3 = 2$ , i.e., 010, for a final representation of 110.010.

## 2.11 Concluding remarks

In the context of inverted indexes, we can exploit the fact that inverted lists are sorted – and typically, strictly increasing – to achieve better compression. In particular, given a sequence  $S[1..n]$  of this form, we can transform the sequence into  $S'$  where  $S'[i] = S[i] - S[i - 1]$  for  $i > 1$  and  $S'[1] = S[1]$ . In the literature,  $S'$  is said to be formed by the so-called *gaps* of  $S$  (or *delta-gaps*). Using the codes described in this section on the gaps of  $S$  is a very popular strategy for compressing inverted indexes, with the key requirement of performing a prefix-sum during decoding. Clearly, compressing these gaps is far more convenient than compressing the integers of  $S$  because the gaps are smaller, thus less bits are required for their codes. In this case, the compressibility of the inverted index critically depends on the distribution of the gaps but, fortunately, most of them are small. Fig. 2 shows the distribution of the gaps for three large text collection that we will introduce in Section 6 (see Table 10, at page 28, for their basic statistics). The plot highlights the *skewed* distribution of the gaps: the most frequent integer is a gap of 1 so that, for better visualization, we cut the percentage to 16% but report the actual value in bold. For example, on the ClueWeb09 dataset 50% of the gaps are just 1. The other values have decreasing frequencies. We divide the distribution into buckets of exponential size, namely the buckets  $B = [1, 2, 4, 8, \dots, 8192]$ . In particular, bucket  $B[i]$  comprises all gaps  $g$  such that  $B[i - 1] < g \leq B[i]$ . (The last bucket also comprises all other gaps larger than 8192 – the “tail” of the distribution.)

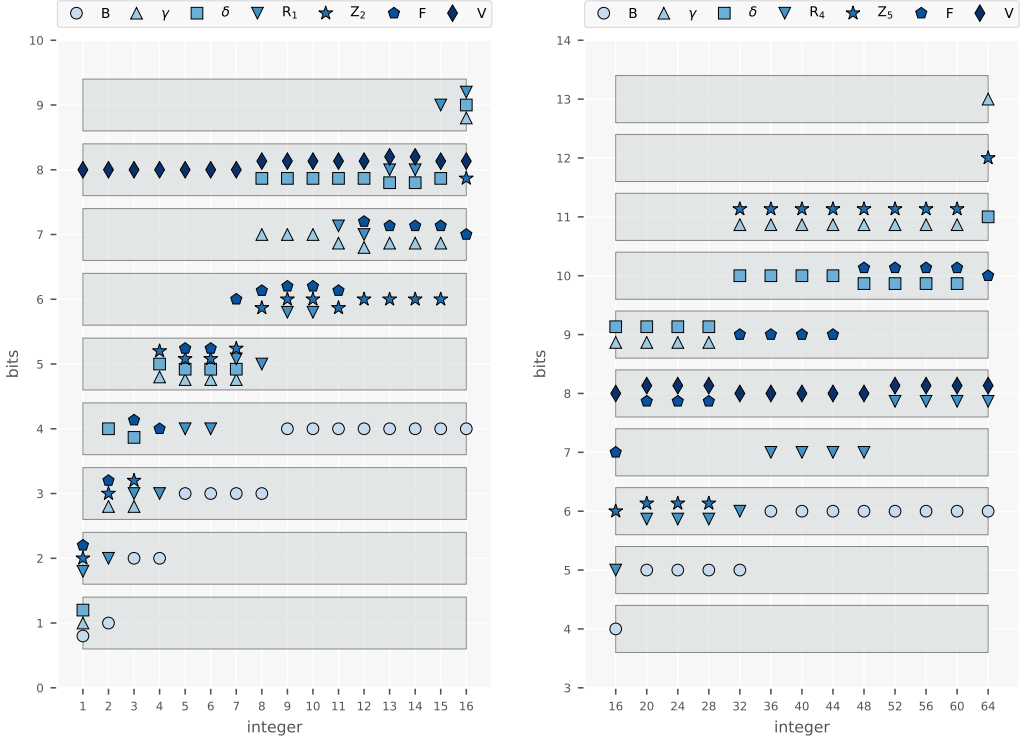


Fig. 3. Comparison between several codes described in Section 2 for the integers 1..64.

As an illustrative comparison between several of the codes described in this section, we report in Fig. 3 the number of bits taken by their codewords when representing the integers 1..64, knowing that such values cover most of the gaps we have in inverted index data (e.g., approximately 86 – 95% of the gaps shown in Fig. 2). In particular, we show the comparison between the codes: binary (B) as an illustrative “single-value” lower bound,  $\gamma$ ,  $\delta$ , Rice (R), Zeta (Z), Fibonacci (F) and Variable-Byte (V). In the plots, data points corresponding to different methods that have the same coordinates have been stacked vertically for better visualization otherwise these would have been indistinguishable. For example,  $\gamma$ ,  $\delta$ , Zeta and Fibonacci all take 5 bits to represent the integers 5 and 6. Not surprisingly, a tuned parametric code such as Rice or Zeta may be the best choice. However, tuning is not always possible and a single parameter has to be specified for the encoding of all integers in the list (or, say, in a sufficiently large block). For the smallest integers the universal codes  $\gamma$  and  $\delta$  are very good, but not competitive immediately for slightly larger integers, e.g., larger than 16. On such values and larger, a simple byte-aligned strategy such as Variable-Byte performs well.

### 3 LIST COMPRESSORS

In this section we describe algorithms that encode an integer list, instead of representing each single integer separately. A useful tool to analyze the space effectiveness of these compressors is the combinatorial *information-theoretic* lower bound, giving the minimum number of bits needed

to represent a list of  $n$  *strictly increasing* integers drawn *at random* from a universe of size  $U \geq n$ , that is [74] ( $e = 2.718$  is the base of the natural logarithm):

$$\left\lceil \log_2 \binom{U}{n} \right\rceil = n \log_2(eU/n) - \Theta(n^2/U) - O(\log n)$$

which is approximately

$$n(\log_2(U/n) + \log_2 e) = n \log_2(U/n) + 1.443n \text{ bits, for } n = o(\sqrt{U}).$$

However, it is important to keep in mind that the compressors we describe in this section often take less space than that computed using the information-theoretic lower bound. In fact, while the lower bound assumes that the integers are distributed at random, the compressors take advantage of the fact that inverted lists feature *clusters* of close integers, e.g., runs of consecutive integers, that are far more compressible than highly scattered regions. This is also the reason why these compressors usually outperform the ones presented in Section 2, at least for sufficiently long lists. As a preliminary example of exploitation of such local clusters, we mention the *Binary Adaptive Sequential Coding* (BASC) by Moffat and Anh [61]. Given a sequence of integers  $\mathcal{S}[1..n]$ , instead of coding the binary magnitude  $b_i = \lceil \log_2(\mathcal{S}[i] + 1) \rceil$  of every single integer  $\mathcal{S}[i]$  – as it happens, for example, in the Elias’ and related codes – we can assume  $b_{i+1}$  to be similar to  $b_i$  (if not the same). This allows to code  $b_{i+1}$  *relatively* to  $b_i$ , hence amortizing its cost.

Such natural clusters of integers are present because the indexed documents themselves tend to be clustered, i.e., there are subsets of documents sharing the very same set of terms. Consider all the Web pages belonging to a certain domain: since their topic is likely to be the same, they are also likely to share a lot of terms. Therefore, not surprisingly, list compressors greatly benefit from docID-reordering strategies that focus on re-assigning the docIDs in order to form larger clusters. When the indexed documents are Web pages, a simple and effective strategy is to assign identifiers to documents according to the lexicographic order of their URLs [95]. This is the strategy we use in the experimental analysis in Section 6 and its benefit is highlighted by Fig. 2: the most frequent gap size is just 1. Another approach uses a recursive graph bisection algorithm to find a suitable re-ordering of docIDs [26]. In this model, the input graph is a bipartite graph in which one set of vertices represents the terms of the index and the other set represents the docIDs. A graph bisection identifies a permutation of the docIDs and, thus, the goal is that of finding, at each step of recursion, the bisection of the graph which minimizes the size of the graph compressed using delta encoding. (There are also other reordering strategies that may be relevant [7, 94]: the works cited here are not meant to be part of an exhaustive list.)

### 3.1 Binary packing

A simple way to improve both compression ratio and decoding speed is to encode a *block* of integers, instead of the whole sequence. This line of work finds its origin in the so-called *frame-of-reference* (FOR) [39]. Once the sequence has been partitioned into blocks (of fixed or variable length), then each block is encoded separately. The key insight behind this simple idea is that, if the sequence is locally homogeneous, i.e., it features clusters of close integers, the values in a block are likely to be of similar magnitude. Vice versa, it is generally hard to expect a long sequence to be globally homogeneous and this is why compression on a per-block basis gives usually better results.

An example of this approach is *binary packing*. Given a block, we can compute the bit width  $b = \lceil \log_2(\max + 1) \rceil$  of the *max* element in the block and then represent all integers in the block using  $b$ -bit codewords. Clearly the bit width  $b$  must be stored prior to the representation of the block. Moreover, the gaps between the integers can be computed to lower the value of  $b$ . Many variants of this simple approach have been proposed [25, 50, 96]. For example, in the *Recursive Bottom-Up*



Table 6. The 9 different ways of packing integers in a 28-bit segment as used by Simple9.

4-bit selector	integers	bits per integer	wasted bits
0000	28	1	0
0001	14	2	0
0010	9	3	1
0011	7	4	0
0100	5	5	3
0101	4	7	0
0110	3	9	1
0111	2	14	0
1000	1	28	0

*Coding* (RBU) code proposed by Moffat and Anh [60], blocks of fixed size are considered, the bit width  $b_i$  of each block determined, and the  $i$ -th block represented via  $b_i$ -bit codewords. The sequence of bit widths  $\{b_i\}_i$  needs to be represented as well and the procedure sketched above is applied recursively to it.

Dividing a list into fixed-size blocks may be suboptimal because regions of close identifiers may be contained in a block containing a much larger value. Thus, it would be preferable to split the list into *variable*-size blocks in order to better adapt to the distribution of the integers in the list. Silvestri and Venturini [96] present an algorithm that finds the optimal splitting of a list of size  $n$  in time  $O(kn)$ , where  $k$  is the maximum block size allowed, in order to minimize the overall encoding cost. Lemire and Boytsov [50] report that the fastest implementation of this approach – named *Vector of Splits Encoding* (VSE) – is that using splits of size 1..14, 16 and 32 integers.

Lemire and Boytsov [50] propose word-aligned versions of binary packing for fast decoding. In the scheme called BP32, 4 groups of 32 bit-packed integers each are stored together in a meta block. Each meta block is aligned to 32-bit boundaries and a 32-bit word is used as descriptor of the meta block. The descriptor stores the 4 bit widths of the 4 blocks in the meta block (8-bit width for each block). The variant called SIMD-BP128 combines 16 blocks of 128 integers each that are aligned to 128-bit boundaries. The use of SIMD instructions provides fast decoding speed.

### 3.2 Simple

Rather than splitting the sequence into blocks of integers as in binary packing, we can split the sequence into fixed-memory units and ask how many integers can be packed in a unit. This is the key idea of the *Simple* family of encoders introduced by Anh and Moffat [3]: pack as many integers as possible in a memory word, i.e., 32 or 64 bits. This approach typically provides good compression and high decompression speed.

For example, Simple9 [3] (sometimes also referred to as Simple4b [4]) adopts 32-bit memory words. It dedicates 4 bits to the *selector* code and 28 bits for data. The selector provides information about how many elements are packed in the data segment using equally-sized codewords. A selector 0000 may correspond to 28 1-bit integers; 0001 to 14 2-bit integers; 0010 to 9 3-bit integers (1 bit unused), and so on, as we can see in Table 6. The four bits distinguish between 9 possible configurations. Similarly, Simple16 [112] has 16 possible configurations using 32-bit words. Simple8b [4], instead, uses 64-bit words with 4-bit selectors. Dedicating 60 bits for data offers 14 different combinations rather than just 9, with only 2 configurations having wasted bits rather than 3.

Anh and Moffat [3] also describe two variations of the Simple9 mechanism, named Relative10 and Carryover12. The idea behind Relative10 is to just use 2 bits for the selector, thus allowing 10 packing configurations with 30 bits. In order to make use of more than 4 options, the selector code

is combined with the one of the previous word, hence enabling the whole range of 10 possibilities. However, when packing  $7 \times 4$ -bit integers or  $4 \times 7$ -bit integers, two bits per word are wasted (only 28 out of the 30 available bits are used). Therefore, in the Carryover12 approach these two bits are used to define the selector code of the following word configuration that makes use of the full 32 bits for packing the integers.

A similar approach to that of the Simple family is used in the QMX mechanism, introduced by Trotman [103]. Considering memory words larger than 64 bits is a popular strategy for exploiting the parallelism of SIMD instructions. QMX packs as many integers as possible into 128- or 256-bit words (Quantities) and stores the selectors (eXtractors) separately in a different stream. The selectors are compressed (Multipliers) with *run-length encoding*, that is with a stream of pairs (*value*, *length*). For example, given the sequence [12, 12, 12, 5, 7, 7, 7, 7, 9], its corresponding RLE representation is [(12, 3), (5, 1), (7, 4), (9, 2)].

### 3.3 PForDelta

The biggest limitation of block-based strategies is their space-inefficiency whenever a block contains just one large value, because this forces the compressor to use a universe of representation as large as that value. This is the main motivation for the introduction of a “patched” frame-of-reference or *PForDelta* (PFor), proposed by Zukowski et al. [114]. The idea is to choose a value  $k$  for the universe of representation of the block, such that a large fraction, e.g., 90%, of its integers can be represented using  $k$  bits per integer. All integers that do not fit in  $k$  bits, are treated as *exceptions* and encoded in a separate array using another compressor, e.g., Variable-Byte or Simple. This strategy is called *patching*. More precisely, two configurable parameters are chosen: a base value  $b$  and a universe of representation  $k$ , so that most of the values fall in the range  $[b, b + 2^k - 1]$  and can be encoded with  $k$  bits each by shifting them (delta-encoding) in the range  $[0, 2^k - 1]$ . To mark the presence of an exception, we also need a special *escape* symbol, thus we have  $[0, 2^k - 2]$  available configurations.

For example, the sequence [3, 4, 7, 21, 9, 12, 5, 16, 6, 2, 34] is represented using PForDelta with parameters  $b = 2$  and  $k = 4$  as [1, 2, 5, \*, 7, 10, 3, \*, 4, 0, \*] – [21, 16, 34]. The special symbol \* marks the presence of an exception that is written in a separate sequence, here reported after the dash.

The *optimized* variant Opt-PFor devised by Yan, Ding, and Suel [111], which selects for each block the values of  $b$  and  $k$  that minimize its space occupancy, it is more space-efficient and only slightly slower than the original PFor. Lemire and Boytsov [50] proposed another variant called Fast-PFor where exceptions are compressed in *pages*, i.e., groups of blocks of integers. For example, a page may be 32 consecutive blocks of 128 integers each, for a total of 4096 integers. In this scheme, all the  $b$ -bit exceptions from all the blocks in a page are stored contiguously, for  $b = 1..32$ . What makes this organization faster to decode is the fact that exceptions are decoded in bulk at a page level, rather than at a (smaller) block level as in Opt-PFor.

In the *parallel* PFor method, proposed by Ao et al. [5], exceptions are represented in a different way to allow their decompression in parallel with that of the “regular” values. Instead of using the escape symbol, each time an exception  $x$  is encountered only the least  $k$  bits of  $x - b$  are written and the overflow bits accumulated in a separate array. The positions of the exceptions are stored in another array and compressed using a suitable mechanism. The same sequence used in the above example, for  $b = 2$  and  $k = 4$ , becomes [1, 2, 5, 3, 7, 10, 3, **14**, 4, 0, **0**] – [1, 0, 2] – [4, 8, 11], because: the least 4 bits of the exceptions  $21 - 2 = 19$ ,  $16 - 2 = 14$ , and  $34 - 2 = 32$  are 3, 14 and 0 respectively (in bold font); the corresponding overflow bits are 1, 0, and 2; the three exceptions appear at positions 4, 8, and 11.

Table 7. An example of Elias-Fano encoding applied to the sequence  $S = [3, 4, 7, 13, 14, 15, 21, 25, 36, 38, 54, 62]$ .

$S$	3	4	7	13	14	15	21	25	36	38		54	62
<i>high</i>	0	0	0	0	0	0	0	0	1	1	1	1	1
	0	0	0	0	0	0	1	1	0	0	0	1	1
	0	0	0	1	1	1	0	1	0	0	1	0	1
<i>low</i>	0	1	1	1	1	1	1	0	1	1		1	1
	1	0	1	0	1	1	0	0	0	1		1	1
	1	0	1	1	0	1	1	1	0	0		0	0
$H$	1110			1110			10	10	110		0	10	10
$L$	011.100.111			101.110.111			101	001	100.110		110		110

### 3.4 Elias-Fano

The encoder we now describe was independently proposed by Elias [30] and Fano [33]. Let  $S(n, U)$  indicate a sorted sequence  $S[1..n]$  whose integers are drawn from a universe of size  $U > S[n]$ . The binary representation of each integer  $S[i]$  as  $\text{bin}(S[i], \lceil \log_2 U \rceil)$  is split into two parts: a *low* part consisting in the right-most  $\ell = \lceil \log_2(U/n) \rceil$  bits that we call *low bits* and a *high* part consisting in the remaining  $\lceil \log_2 U \rceil - \ell$  bits that we similarly call *high bits*. Let us call  $\ell_i$  and  $h_i$  the values of low and high bits of  $S[i]$  respectively. The Elias-Fano encoding of  $S(n, U)$  is given by the encoding of the high and low parts. The integers  $[\ell_1, \dots, \ell_n]$  are written verbatim in a bitvector  $L$  of  $n \lceil \log_2(U/n) \rceil$  bits, which represents the encoding of the low parts. The high parts are represented with another bitvector of  $n + 2^{\lceil \log_2 n \rceil} \leq 2n$  bits as follows. We start from a 0-valued bitvector  $H$  and set the bit in position  $h_i + i$ , for all  $i = 1, \dots, n$ . It is easy to see that the  $k$ -th unary value  $m$  of  $H$  indicates that  $m - 1$  integers of  $S$  have high bits equal to  $k$ ,  $0 \leq k \leq \lceil \log_2 n \rceil$ . Finally the Elias-Fano representation of  $S$  is given by the concatenation of  $H$  and  $L$  and overall takes

$$\text{EF}(S(n, U)) \leq n \lceil \log_2(U/n) \rceil + 2n \text{ bits.} \quad (1)$$

Although we can opt for an arbitrary split into high and low parts, ranging from 0 to  $\lceil \log_2 U \rceil$ , it can be shown that  $\ell = \lceil \log_2(U/n) \rceil$  minimizes the overall space occupancy of the encoding [30]. Moreover, given that the information-theoretic lower bound is approximately  $n \log_2(U/n) + n \log_2 e$  bits, it can be shown [30] that less than half a bit is wasted per element by Formula 1. Table 7 shows an example of encoding for the sequence  $[3, 4, 7, 13, 14, 15, 21, 25, 36, 38, 54, 62]$ . Note that no integer has high part equal to 101.

The same code arrangement was later described by Anh and Moffat [2] as a “modified” version of the Rice code (Section 2.5). In fact, they partition  $U$  into buckets of  $2^k$  integers each for some  $k > 0$ , code in unary how many integers fall in each bucket, and represent each integer using  $k$  bits as an offset to its bucket. The connection with Rice is established by writing the number of integers sharing the same quotient, rather than encoding this quantity for every integer. They also indicated that the optimal parameter  $k$  should be chosen to be  $\lceil \log_2(U/n) \rceil$ .

**Supporting random Access.** Despite the elegance of the encoding, it is possible to support random access to individual integers *without* decompressing the whole sequence. Formally, we are interested in implementing the operation  $\text{Access}(i)$  that returns  $S[i]$ . The operation can be implemented by using an auxiliary data structure that is built on the bitvector  $H$  and efficiently answers  $\text{Select}_1$  queries. The answer to a  $\text{Select}_b(i)$  query over a bitvector is the position of the  $i$ -th bit set to  $b$ . For example,  $\text{Select}_0(3) = 10$  on the bitvector  $H$  of Table 7. This auxiliary data structure is *succinct* in the sense that it is negligibly small in asymptotic terms, compared to  $\text{EF}(S(n, U))$ ,

requiring only  $o(n)$  additional bits [54, 107]. Using the  $\text{Select}_1$  primitive, it is possible to implement Access in  $O(1)$ . (A prior method than that using Select is described by Anh and Moffat [2], who adopted a byte-wise processing algorithm to accelerate skipping thorough the  $H$  bitvector.)

We basically have to “re-link” together the high and low bits of an integer, previously split up during the encoding phase. The low bits  $\ell_i$  are trivial to retrieve as we need to read the range of bits  $\ell_i = L[(i-1)\ell + 1, i\ell]$ . The retrieval of the high bits is, instead, more complicated. Since we write in unary how many integers share the same high part, we have a bit set for every integer in  $\mathcal{S}$  and a zero for every distinct high part. Therefore, to retrieve the high bits of the  $i$ -th integer, we need to know how many zeros are present in the first  $\text{Select}_1(i)$  bits of  $H$ . This quantity is evaluated on  $H$  in  $O(1)$  as  $\text{Select}_1(i) - i$ . Linking the high and low bits is as simple as:  $\text{Access}(i) = ((\text{Select}_1(i) - i) \ll \ell) \mid \ell_i$ , where  $\ll$  indicates the left shift operator and  $\mid$  is the bitwise OR.

For example, to recover  $\mathcal{S}[4] = 13$ , we first evaluate  $\text{Select}_1(4) - 4 = 5 - 4 = 1$  and conclude that the high part of  $\mathcal{S}[4]$  is the binary representation of 1, that is 001. Finally, we access the low bits  $L[10..12] = 101$  and re-link the two parts, hence obtaining 001.101.

**Supporting Successor queries.** The query  $\text{Successor}(x)$ , returning the smallest integer  $y$  of  $\mathcal{S}$  such that  $y \geq x$ , is supported in  $O(1 + \log(U/n))$  time as follows. Let  $h_x$  be the high bits of  $x$ . Then for  $h_x > 0$ ,  $i = \text{Select}_0(h_x) - h_x + 1$  indicates that there are  $i$  integers in  $\mathcal{S}$  whose high bits are less than  $h_x$ . On the other hand,  $j = \text{Select}_0(h_x + 1) - h_x$  gives us the position at which the elements having high bits greater than  $h_x$  start. The corner case  $h_x = 0$  is handled by setting  $i = 0$ . These two preliminary operations take  $O(1)$ . Now we can conclude the search in the range  $\mathcal{S}[i, j]$ , having *skipped* a potentially large range of elements that, otherwise, would have required to be compared with  $x$ . We therefore determine the successor of  $x$  by binary searching in this range which contains up to  $U/n$  integers. The time bound follows.

As an example, consider the query  $\text{Successor}(30)$  over the example sequence from Table 7. Since  $h_{30} = 3$ , we have  $i = \text{Select}_0(3) - 3 + 1 = 8$  and  $j = \text{Select}_0(4) - 3 = 9$ . Therefore we conclude our search in the range  $\mathcal{S}[8, 9]$  by returning  $\text{Successor}(30) = \mathcal{S}[9] = 36$ .

In the specific context of inverted indexes, the query Successor is called NextGEQ (Next Greater-than or Equal-to) and we are going to adopt this terminology in Section 6. It should also be observed that Moffat and Zobel [70] were the first to explore the use of *skip pointers* – meta data aimed at accelerating the skipping through blocks of compressed integers – for faster query evaluation.

**Partitioning the integers by cardinality.** One of the most pertinent characteristics of the Elias-Fano space bound in Formula 1 is that it only depends on two parameters, i.e., the size  $n$  of the sequence  $\mathcal{S}$  and the universe  $U > \mathcal{S}[n]$ . As already explained, inverted lists often present clusters of very similar integers and Elias-Fano fails to exploit them for better compression because it always uses a number of bits per integer at most equal to  $\lceil \log_2(U/n) \rceil + 2$ , thus proportional to the logarithm of the *average gap*  $U/n$  between the integers and *regardless* any skewed distribution. (Note that also Golomb and Rice are insensitive to any deviation away from a random selection of the integers.) In order to better adapt to the distribution of the gaps between the integers, we can partition the sequence, obtaining the so-called *partitioned Elias-Fano* (PEF) representation [73].

The sequence is split into  $k$  blocks of variable length. The first level of representation stores two sequences compressed with plain Elias-Fano: (1) the sequence made up of the last elements  $\{U_1, \dots, U_k\}$  of the blocks, the so-called *upper-bounds* and (2) the prefix-summed sequence of the sizes of the blocks. The second level is formed, instead, by the representation of the blocks themselves, that can be again encoded with Elias-Fano. The main advantage of this two-level representation, is that now the integers in the  $i$ -th block are encoded with a smaller universe, i.e.,  $U_i - U_{i-1}$ ,  $i > 0$ , thus improving the space with respect to the original Elias-Fano representation. More precisely, each block in the second level is encoded with one among *three* different strategies.

As already stated, one of them is Elias-Fano. The other two additional strategies come into play to overcome the space inefficiencies of Elias-Fano when representing *dense* blocks.

Let consider a block and call  $b$  its size,  $M$  its universe respectively. Vigna [107] observed that as  $b$  approaches  $M$  the space bound  $b\lceil\log_2(M/b)\rceil + 2b$  bits becomes close to  $2M$  bits. In other words, the closer  $b$  is to  $M$ , the denser the block. However, we can always represent the block with  $M$  bits by writing the *characteristic vector* of the block, that is a bitvector where the  $i$ -th bit is set if the integer  $i$  belongs to the block. Therefore, besides Elias-Fano, two additional encodings can be chosen to encode the block, according on the relation between  $m$  and  $b$ . The first one addresses the extreme case in which the block covers the whole universe, i.e., when  $b = M$ : in such case, the first level of the representation (upper-bound and size of the block) trivially suffices to recover each element of the block that, therefore, does not need to be represented at all. The second case is used whenever the number of bits used by the Elias-Fano representation of the block is larger than  $M$  bits: by doing the math, it is not difficult to see that this happens whenever  $b > M/4$ . In this case we can encode the block with its characteristic bitvector using  $M$  bits. The choice of the proper encoding for a block is rather fundamental for the practical space effectiveness of PEF.

Let us consider a simple example with  $M = 40$  bits. Suppose that the block is sparse, e.g., with  $b = 5$ . Then, Elias-Fano takes  $\lceil\log_2(40/5)\rceil + 2 = 5$  bits per element, whereas a characteristic vector representation would take  $40/5 = 8$  bits per element. In a dense case with, say,  $b = 30$ , a bitmap just takes  $40/30 = 1.33$  bits per element, whereas Elias-Fano would take 3 bits per element.

Splitting the sequence into equally-sized block is clearly sub-optimal, since we cannot expect clusters of similar integers to be aligned with uniform partitions. For such reason, an algorithm based on dynamic programming is presented by Ottaviano and Venturini [73] that yields a partition whose cost in bits is at most  $(1 + \epsilon)$  times away from the optimal one taking  $O(n \log \frac{1}{\epsilon})$  time and  $O(n)$  space for any  $0 < \epsilon < 1$ . Notice that the time complexity becomes  $\Theta(n)$  when  $\epsilon$  is constant. In fact, the problem of determining the partition of minimum encoding cost can be seen as the problem of finding the path of minimum cost (shortest) in a complete, weighted and directed acyclic graph (DAG). This DAG has  $n$  vertices, one for each integer of  $\mathcal{S}$ , and  $\Theta(n^2)$  edges where the cost  $w(i, j)$  of edge  $(i, j)$  represents the number of bits needed to represent  $\mathcal{S}[i, j]$ . Each edge cost  $w(i, j)$  is computed in  $O(1)$  by just knowing the universe and size of the chunk  $\mathcal{S}[i, j]$ . By pruning the DAG it is possible to attain to the mentioned complexity by preserving the approximation guarantees [73].

**Partitioning the integers by universe.** As already mentioned, we can opt for an arbitrary split between the high and the low part of the Elias-Fano representation. Partitioning the universe  $U$  into chunks containing at most  $2^\ell$  integers each, with  $\ell = \lceil\log_2(U/n)\rceil$ , minimizes the space of the encoding [30] but a *non-parametric* split – independent from the values of  $U$  and  $n$  – is also possible. Let us assume that  $U \leq 2^{32}$  in the following.

For example, *Roaring* [15, 51, 53] partitions  $U$  into chunks spanning  $2^{16}$  values each and represents all the integers of the sequence falling into a chunk in two different ways according to the cardinality of the chunk: if the chunk contains less than 4096 elements, then it is considered to be *sparse* and represented as a sorted array of 16-bit integers; otherwise it is considered *dense* and encoded as a bitmap of  $2^{16}$  bits. Lastly, very dense chunks can also be encoded with *runs* if advantageous. A run is represented as a pair  $(v, \ell)$  meaning that all the integers  $v \leq x \leq v + \ell$  belong to the chunk.

Inspired by the van Emde Boas tree [105, 106], the *Slicing* [78] data structure recursively slices the universe of representation in order to better adapt to the distribution of the integers being compressed. Differently from *Roaring*, a *sparse* sparse chunk is further partitioned into at most  $2^8$  blocks of  $2^8$  elements each. Therefore, a non-empty universe slice of  $2^{16}$  elements can be either: represented with a bitmap of  $2^{16}$  bits (dense case); represented implicitly if the slice contains all the

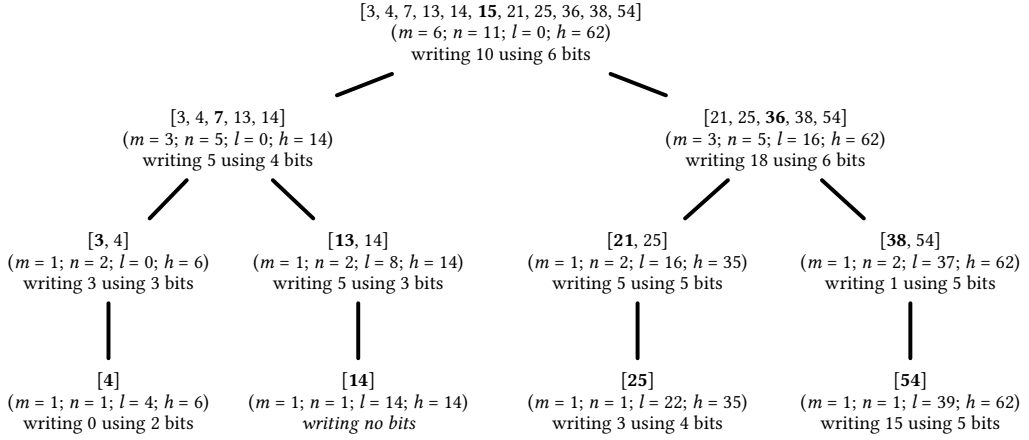


Fig. 4. The recursive calls performed by the Binary Interpolative Coding algorithm when applied to the sequence [3, 4, 7, 13, 14, 15, 21, 25, 36, 38, 54] with initial knowledge of lower and upper bound values  $l = 0$  and  $h = 62$ . In bold font we highlight the middle element being encoded.

possible  $2^{16}$  elements (full case); or it is recursively partitioned into smaller slices of  $2^8$  elements each. Finally, each non-empty slice of  $2^8$  elements is encoded with a sorted array of 8-bit integers (sparse case); or with a bitmap of  $2^8$  bits (dense case). The idea of a hybrid compression scheme with hierarchical bit-vectors and sorted-arrays (that can be further compressed) was first proposed by Fraenkel et al. [35].

It should be noted that all the partitioning strategies we have described in this section, namely partitioned Elias-Fano (PEF), Roaring and Slicing, exploit the same idea to attain to good space effectiveness: look for dense regions to be encoded with bitmaps and use a different mechanism for sparse regions. While PEF achieves this goal by splitting the sequence *optimally* by cardinality, Roaring and Slicing partition the universe of representation *greedily*, hence maintaining the property that all partitions are represented using the *same* universe. As we will better see in Section 6, these different partitioning paradigms achieve different space/time trade-offs.

### 3.5 Interpolative

The *Binary Interpolative Code* (BIC) invented by Moffat and Stuiver [65, 66] represents a sorted integer sequence without requiring the computation of its gaps. The key idea of the algorithm is to exploit the order of the already-encoded elements to compute the number of bits needed to represent the elements that will be encoded next.

At the beginning of the encoding phase, suppose we are specified two quantities  $l \leq S[1]$  and  $h \geq S[n]$ . Given such quantities, we can encode the element in the middle of the sequence, i.e.,  $S[m]$  with  $m = \lceil n/2 \rceil$ , in some appropriate manner, knowing that  $l \leq S[m] \leq h$ . For example, we can write  $S[m] - l - m + 1$  using just  $\lceil \log_2(h - l - n + 1) \rceil$  bits. After that, we can apply the same step to both halves  $S[1, m-1]$  and  $S[m+1, n]$  with *updated knowledge* of lower and upper values  $(l, h)$  that are set to  $(l, S[m] - 1)$  and  $(S[m] + 1, h)$  for the left and right half respectively. Note that whenever the condition  $l + n - 1 = h$  is satisfied, a “run” of consecutive integers is detected: therefore, we can stop recursion and emit no bits at all during the encoding phase. When the condition is met again during the decoding phase, we simply output the values  $l, l + 1, l + 2, \dots, h$ . This means that BIC can actually use codewords of 0 bits to represent more than one integer, hence



attaining to a rate of less than one bit per integer – a remarkable property that makes the code very succinct for highly clustered inverted lists.

We now consider an encoding example applied to the sequence [3, 4, 7, 13, 14, 15, 21, 25, 36, 38, 54, 62]. As it is always safe to choose  $l = 0$  and  $h = S[n]$ , we do so, thus at the beginning of the encoding phase we have  $l = 0$  and  $h = 62$ . Since we set  $h = S[n]$ , the last value of the sequence is first encoded and we process  $S[1..n-1]$  only. Fig. 4 shows the sequence of recursive calls performed by the encoding algorithm oriented as a binary tree. At each node of the tree we report the values assumed by the quantities  $m$ ,  $n$ ,  $l$  and  $h$ , plus the processed subsequence and the number of bits needed to encode the middle element. By *pre-order* visiting the tree, we obtain the sequence of written values, that is [10, 5, 3, 0, 5, 18, 5, 3, 1, 15] with associated codeword lengths [6, 4, 3, 2, 3, 6, 5, 4, 5, 5]. Note that the value in the second leaf of the tree, i.e.,  $S[5] = 14$ , is encoded with 0 bits given that both  $l$  and  $h$  are equal to 14.

However, the encoding process obtained by the use of simple binary codes as illustrated in Fig. 4 is wasteful. In fact, as discussed in the original work [65, 66], more succinct encodings can be achieved with a minimal binary encoding (recall the definition at the end of Section 2.2). More precisely, when the range  $r > 0$  is specified, all values  $0 \leq x \leq r$  are assigned fixed-length codewords of size  $\lceil \log_2(r+1) \rceil$  bits. But the more  $r$  is distant from  $2^{\lceil \log_2(r+1) \rceil}$  the more this allocation of codewords is wasteful because  $c = 2^{\lceil \log_2(r+1) \rceil} - r - 1$  codewords can be made 1 bit shorter without loss of unique decodability. Therefore we proceed as follows. We identify the range of smaller codewords, delimited by the values  $r_l$  and  $r_h$ , such that every value  $x \leq r$  such that  $r_l < x < r_h$  is assigned a shorter  $(\lceil \log_2(r+1) \rceil - 1)$ -bit codeword and every value outside this range is assigned a longer one of  $\lceil \log_2(r+1) \rceil$  bits. To maintain unique decodability, we first always read  $\lceil \log_2(r+1) \rceil - 1$  bits and interpret these as the value  $x$ . Then we check if condition  $r_l < x < r_h$  is satisfied: if so, we are done; otherwise, the codeword must be extended by 1 bit. In fact, in a *left-most* minimal binary code assignment, the first  $c$  values are assigned the shorter codewords, thus  $r_h = 2^{\lceil \log_2(r+1) \rceil} - r - 1$  (and we only check whether  $x < r_h$ ). In a *centered* minimal binary code assignment, the values in the centre of the range are assigned the shorter codewords, thus  $(r_l, r_h) = (\lfloor r/2 \rfloor - \lfloor c/2 \rfloor - 1, \lfloor r/2 \rfloor + \lfloor c/2 \rfloor + 1)$  if  $r$  is even, or  $(r_l, r_h) = (\lfloor r/2 \rfloor - \lfloor c/2 \rfloor, \lfloor r/2 \rfloor + \lfloor c/2 \rfloor + 1)$  if  $r$  is odd. The rationale behind using centered minimal codes is that a reasonable guess is to assume the middle element to be about half of the upper bound. As already noted, we remark that the specific assignment of codewords is irrelevant and many assignments are possible: what matters is to assign correct lengths and maintain the property of unique decodability.

It is also worth mentioning that the *Tournament code* developed by Teuhola [100] is very related to BIC.

### 3.6 Directly-addressable codes

Brisaboa, Ladra, and Navarro [12] introduced a representation for a list of integers that supports random access to individual integers – called *directly-addressable code* (DAC) – noting that this is not generally possible for many of the representations described in Section 2 and 3. They reduced the problem of random access to the one of *ranking* over a bitmap. Given a bitmap  $B[1..n]$  of  $n$  bits, the query  $\text{Rank}_b(B, i)$  returns the number of  $b$  bits in  $B[1, i]$ , for  $i \leq n$ . For example, if  $B = 010001101110$  then  $\text{Rank}_1(6) = 2$  and  $\text{Rank}_0(8) = 5$ . Rank queries can be supported in  $O(1)$  by requiring only  $o(n)$  additional bits [17, 41, 44].

Each integer in the list is partitioned into  $(b+1)$ -bit chunks. Similarly to Variable-Byte,  $b$  bits are dedicated to the representation of the integer and the control bit indicates whether another chunk follows or not. All the first  $n$   $b$ -bit chunks of every integer are grouped together in a codeword stream  $C_1$  and the control bits form a bitmap  $B_1[1..n]$  of  $n$  bits. If the  $i$ -th bit is set in such bitmap,

then the  $i$ -th integer in the sequence needs a second chunk, otherwise a single chunk is sufficient. Proceeding recursively, all the second  $m \leq n$  chunks are concatenated together in  $C_2$  and the control bits in a bitmap  $B_2[1..m]$  of  $m$  bits. Again, the  $i$ -th bit of such bitmap is set if the  $i$ -th integer with at least two chunks needs a third chunk of representation. In general, if  $U$  is the maximum integer in the list, there are at most  $\lceil \log_2(U+1)/b \rceil$  levels (i.e., streams of chunks).

As an example, the sequence  $[2, 7, 12, 5, 13, 142, 61, 129]$  is encoded with  $b = 3$  as follows:  $B_1 = 0.0.1.0.1.1.1.1$ ;  $B_2 = 0.0.1.0.1$ ;  $B_3 = 0.0$ ;  $C_1 = 010.111.100.101.101.110.011.001$ ;  $C_2 = 001.001.001.110.000$ ;  $C_3 = 010.010$ .

Accessing the integer  $x$  in position  $i$  reduces to a sequence of  $c - 1$   $\text{Rank}_1$  operations over the levels' bitmaps, where  $c \geq 1$  is the number of  $(b + 1)$ -bit chunks of  $x$ , that is  $\lceil \log_2(x + 1)/b \rceil$ . Now, for  $k = 1..c$ , we repeat the following step: (1) retrieve the  $i$ -th chunk from the  $C_k$  in constant time given that all chunks are  $b$  bits long; (2) if  $B_k[i] = 0$ , we are done; otherwise  $j = \text{Rank}_1(B_k, i)$  gives us the number of integers (in the level  $k + 1$ ) that have more than  $k$  chunks, so we set  $i = j$  and repeat. For example,  $\text{Access}(5)$  is resolved as follows on our example sequence. We retrieve  $C_1[5] = 101$ ; since  $B_1[5] = 1$ , we compute  $\text{Rank}_1(B_1, 5) = 2$ . Now we retrieve  $C_2[2] = 001$  and given  $B_2[2] = 0$ , we stop by returning the integer  $C_2[2].C_1[5] = 001.101$ , that is 13.

Lastly, nothing prevents from changing the value of  $b$  at each level of the data structure. For this purpose, the authors of DAC present an algorithm, based on dynamic programming, that finds such optimal values for a given list.

### 3.7 Hybrid approaches

Hybrid approaches are possible by using different compressors to represent the blocks of a list. For example, given a query log, we can collect access statistics at a block-level granularity, namely how many times a block is accessed during query processing, and represent rarely-accessed blocks with more space-efficient compressor; vice versa frequently-accessed blocks are encoded with more time-efficient compressor [72]. This hybrid strategy produces good space/time trade-offs.

Pibiri and Venturini [83] show that a list of  $n$  sorted integers can be optimally partitioned into variable-length blocks whenever the chosen representation for each block is given by either: (1) any compressor described in Section 2, namely a *point-wise* encoder, or (2) the characteristic vector of the block. From Section 3.4 we recall that, given a block of universe  $m$ , the characteristic vector representation of the block is given by a bitmap of  $m$  bits where the  $i$ -th bit is set if the integer  $i$  belongs to the block. By exploiting the fact that the chosen encoder is point-wise, i.e., the number of bits needed to represent an integer solely depends on the integer itself rather than the block where it belongs to, it is possible to devise an algorithm that finds an *optimal* partitioning in  $\Theta(n)$  time and  $O(1)$  space. The constant factor hidden by the asymptotic notation is very small, making the algorithm very fast in practice.

### 3.8 Entropy coding: Huffman, Arithmetic, and Asymmetric Numeral Systems

In this section we quickly survey the most famous *entropy coding* techniques – Huffman [43], Arithmetic coding [62, 76, 88, 89], and Asymmetric Numeral Systems (ANS) [27–29]. Although some authors explored the use of these techniques for index compression, especially Huffman [10, 36, 45, 69] and ANS [63, 64] (see Section 4.2), they are usually not competitive in terms of efficiency and implementation simplicity against the compressors we have illustrated in the previous sections, making them a hard choice for practitioners. An in-depth treatment of such techniques is, therefore, outside the scope of this article and the interested reader can follow the references to the individual papers we include here. The survey by Moffat [59] about Huffman coding also contains descriptions of Arithmetic Coding and ANS (Section 5.1 and 5.2 of that article, respectively).

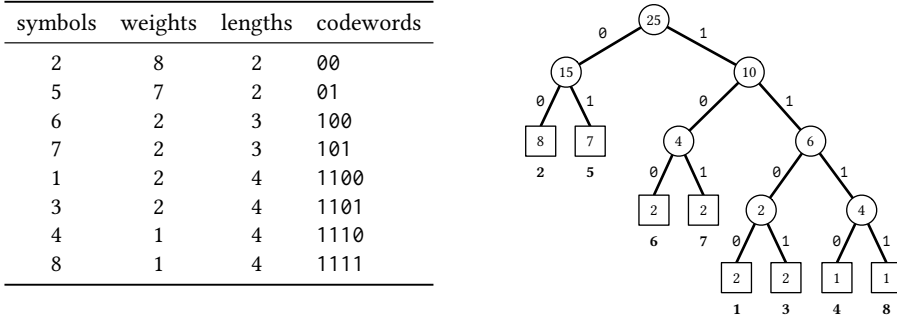


Fig. 5. An example of Huffman coding applied to a sequence of size 25 with symbols 1..8 and associated weights [2, 8, 2, 1, 7, 2, 2, 1].

We first recall the definition of *entropy*, a tool introduced by Shannon [93]. He was concerned with the problem of defining the *information content* of a discrete random variable  $X : \Sigma \rightarrow \mathbb{R}$ , with distribution  $\mathbb{P}(s) = \mathbb{P}\{X = s\}$ ,  $s \in \Sigma$ . He defined the entropy of  $X$  as  $H = \sum_{s \in \Sigma} [\mathbb{P}(s) \log_2(1/\mathbb{P}(s))]$  bits. The quantity  $\log_2(1/\mathbb{P}(s))$  bits is also called the *self-information* of the symbol  $s$  and  $H$  represents the average number of bits we need to encode each value of  $\Sigma$ . Let now  $S$  be a sequence of  $n$  symbols drawn from an alphabet  $\Sigma$ . (In the context of this article, the symbols will be integer numbers.) Let also  $n_s$  denote the number of times the symbol  $s$  occurs in  $S$ . Assuming empirical frequencies as probabilities [75] (the larger is  $n$ , the better the approximation), i.e.,  $\mathbb{P}(s) \approx n_s/n$ , we can consider  $S$  as a random variable assuming value  $s$  with probability  $\mathbb{P}(s)$ . In this setting, the entropy of the sequence  $S$  is  $H_0 = 1/n \sum_{s \in \Sigma} [n_s \log_2(n/n_s)]$  bits, also known as the 0-th order (or *empirical*) entropy of  $S$ . In particular, the quantity  $nH_0$  gives a theoretic lower bound on the average number of bits we need to represent  $S$  and, hence, to the output size of *any* compressor that encodes each symbol of  $S$  with a fixed-length codeword.

**Huffman.** It is standard to describe the Huffman's algorithm in terms of a binary tree. In this logical binary tree, a leaf corresponds to a symbol to be encoded with associated symbol occurrence – its weight – and an internal node stores the sum of the weights of its children. The algorithm maintains a *candidate set* of tree nodes from which, at each step: (1) the two nodes with smallest weight are selected; (2) they are merged together into a new parent node whose weight is the sum of the weights of the two children; (3) the parent node is added to the candidate set. The algorithm repeats this merging step until only the root of the tree (whose weight is the length of the sequence  $S$ ) is left in the candidate set. Fig. 5 shows an example of Huffman coding. It is important to mention that, in practice, the decoding process does *not* traverse any tree. An elegant variation of the algorithm – known as *canonical Huffman* – allows fast decoding by using lookup tables as we similarly illustrated in Section 2.1. Again, Moffat [59] provides all details.

Now, let  $L$  be the *average* Huffman codeword length. Two of the most important properties of Huffman coding are: (1)  $L$  is *minimum* among all possible prefix-free codes; (2)  $L$  satisfies  $H_0 \leq L < H_0 + 1$ . The first property means that Huffman coding produces an *optimal* code for a given distribution of the integers. (The precursor of the Huffman's algorithm is the less-known *Shannon-Fano* algorithm that was independently proposed by Shannon [93] and Fano [32], which, however, does not always produce an optimal code.) The second property suggests that an Huffman code can loose up to 1 bit compared to the entropy  $H_0$  because it requires *at least* 1 bit to encode a symbol (as any other prefix-free code), thus if  $H_0$  is large, the extra bit lost is negligible in

practice; otherwise the distribution of probabilities is skewed and Huffman loses a significant space compared to the entropy of the source.

Now, it should be clear why Huffman may *not* be an appropriate choice for inverted index compression. Applying Huffman to the compression of the integers in inverted lists means that its alphabet of representation is too large, thus making the mere description of the code outweigh the cost of representing very sparse inverted lists. The same reason applies if we try to use the code to compress the gaps between successive integers: the largest gap could be as large as the largest integer in the sequence.

**Arithmetic.** The first concept of Arithmetic coding was introduced by Elias before 1963 according to Note 1 on page 61 in the book by Abramson [1]. However, the method requires infinite precision arithmetic and, because of this, it remained unpublished. The first practical implementations were designed during 1976 by Rissanen [89] and Pasco [76], and later refined by Rissanen [88]. A more recent efficient implementation is described by Moffat et al. [62]. The method offers higher compression ratios than Huffman's, especially on highly skewed distributions, because it is *not* a prefix-free code, so it does not require at least one bit to encode a symbol. Indeed a single bit may correspond to more than one input symbol. However, Huffman codes are faster to decode; Arithmetic does not permit to decode an output stream starting from an arbitrary position, but only sequential decoding is possible.

Given a sequence of symbols  $S = [s_1..s_n]$ , the main idea behind the method works as follows. The interval  $[0, 1)$  is partitioned into  $|\Sigma|$  segments of length proportional to the probabilities of the symbols. Then the subinterval corresponding to  $s_1$ , say  $[\ell_1, r_1)$ , is chosen and the same partitioning step is applied to it. The process stops when all input symbols have been processed and outputs a single real number  $x$  in  $[\ell_n, r_n)$ , that is the interval associated to the last input symbol  $s_n$ . Then the pair  $(x, n)$  suffices to decode the original input sequence  $S$ .

It can be shown that Arithmetic coding takes at most  $nH_0 + 2$  bits to encode a sequence  $S$  of length  $n$ . This means that the overhead with respect to the empirical entropy  $H_0$  is only of  $2/n$  bits per symbol, thus negligible for basically all practical values of  $n$ . As already pointed out, Arithmetic coding requires infinite precision that can be very costly to be approximated. In fact, a practical implementation [110] using approximated arithmetic can take up to  $nH_0 + \frac{2}{100}n$  bits, thus having 0.02 bits of loss per symbol rather than  $2/n$ .

**Asymmetric Numeral Systems.** Asymmetric Numeral Systems (ANS) is a family of entropy coding algorithms, originally developed by Duda [27, 28], which approaches the compression ratio of Arithmetic coding with a decompression speed comparable with the one of Huffman [29]. The basic idea of ANS is to represent a sequence of symbols with a natural number  $x$ .

Let us consider a concrete example [64] with an alphabet of 3 symbols only, namely  $\{a, b, c\}$  and assuming that  $\mathbb{P}(a) = 1/2$ ,  $\mathbb{P}(b) = 1/3$  and  $\mathbb{P}(c) = 1/6$ . In order to derive the encoding of a sequence of symbols, a *frame*  $f[1..m]$  of symbols is constructed, having  $1/2$  of the entries equal to  $a$ ,  $1/3$  equal to  $b$  and  $1/6$  equal to  $c$ . For example, one such frame could be  $f[1..6] = [aaabbc]$ , but other symbol permutations with possibly larger  $m$  are possible as well. The frame determines a table that is used to map a sequence of symbols to an entry in the table. Refer to Table 8a for an example with the frame  $f[1..6] = [aaabbc]$ . The entries in the table are the natural numbers assigned incrementally in the order determined by the frame. For example, since the first three symbols in the frame are  $aaa$ , the first numbers assigned to  $a$ 's row are 1, 2 and 3. The next two symbols are  $bb$ , so  $b$ 's row gets 4 and 5. The last symbol in the frame is  $c$ , so the first entry in  $c$ 's row is 6. The process now proceed in cycles, thus placing 7, 8 and 9 in  $a$ 's row; 10 and 11 in  $b$ 's row, a final 12 in  $c$ 's row, and so on (first 10 columns are shown in the table). Table 8b shows an example for another distribution of the symbols, constructed using a frame  $f[1..4] = [caba]$ .

Table 8. Two example of ANS encoding table, with respectively frame  $f[1..6] = [aaabbc]$  and  $f[1..4] = [caba]$ .

(a)												(b)											
$\Sigma$	$\mathbb{P}$	codes										$\Sigma$	$\mathbb{P}$	codes									
$a$	1/2	1	2	3	7	8	9	13	14	15	19	$a$	1/2	2	4	6	8	10	12	14	16	18	20
$b$	1/3	4	5	10	11	16	17	22	23	28	29	$b$	1/4	3	7	11	15	19	23	27	31	35	39
$c$	1/6	6	12	18	24	30	36	42	48	54	60	$c$	1/4	1	5	9	13	17	21	25	29	33	37
		0	1	2	3	4	5	6	7	8	9			0	1	2	3	4	5	6	7	8	9

Now, consider the sequence *caa* and let us determine its ANS code with the table in Fig. 8a. We make use of the *transition function* defined by the table  $T$  itself as  $state' = T[s, state]$  which, given a symbol  $s$  and a  $state$  value, produces the next  $state'$  of the encoder. At the beginning we set  $state = 0$ , thus for the given sequence *caa* the  $state$  variable assumes values  $0 \rightarrow 6 \rightarrow 13 \rightarrow 26$  (last value not shown in the table). The code assigned to the sequence is, therefore, the integer 26. For the sequence *acb* under the encoding table in Table 8b we generate, instead, the transitions  $0 \rightarrow 2 \rightarrow 9 \rightarrow 39$ , thus the assigned code is 39. Decoding reverts this process. For example, given 39 we know that the last symbol of the encoded sequence must have been *b* because 39 is found on the second row of the table. The value is in column 9, which is found in column 2 in the third row that corresponds to the *c* symbol. Finally, the column number 2 is found in *a*'s row, thus we emit the message *acb*.

## 4 INDEX COMPRESSORS

This section is devoted to approaches that look for regularities among *all* the lists in the inverted index. In fact, as already motivated at the beginning of Section 3, the inverted index naturally presents some amount of redundancy in that many sub-sequences of integers are shared between the lists. Good compression can be achieved by exploiting this correlation, usually at the expense of a reduced processing efficiency.

### 4.1 Clustered

Pibiri and Venturini [80] propose a clustered index representation. The inverted lists are grouped into clusters of “similar” lists, i.e., the ones sharing as many integers as possible. Then for each cluster, a *reference list* is synthesized with respect to which all lists in the cluster are encoded. More specifically, the integers belonging to the intersection between the cluster reference list and a list in the cluster are represented as the positions they occupy within the reference list. This makes a big improvement for the cluster space, since each intersection can be re-written in a much smaller universe. Although *any* compressor can be used to represent the intersection and the residual segment of each list, the authors adopt partitioned Elias-Fano; by varying the size of the reference lists, different time/space trade-offs can be obtained.

### 4.2 ANS-based

In Section 3.8 we have seen an example of the ANS method developed by Duda [27, 28]. As we have already observed in that section, the alphabet size may be too large for representing the integers in inverted indexes. Even the largest gap may be equal to the number of documents in the collection, which is usually several order to magnitudes larger than, for example, the (extended) ASCII alphabet. For this reason, Moffat and Petri [63] describe several adaptations of the base ANS mechanism tailored for effective index compression. In order to reduce the alphabet size, they perform a preprocessing step with Variable-Byte to reduce the input list to a sequence of bytes

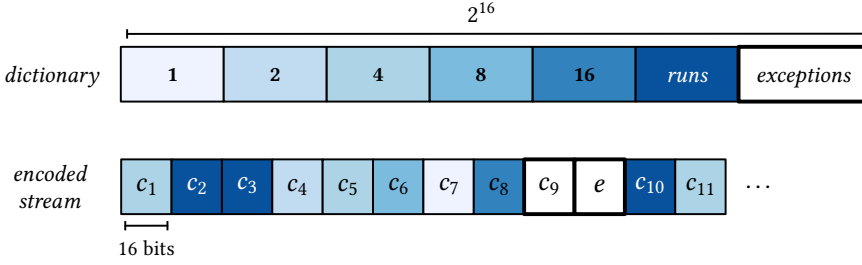


Fig. 6. A dictionary-based encoded stream example, where dictionary entries corresponding to  $\{1, 2, 4, 8, 16\}$ -long integer patterns, runs and exceptions, are labelled with different shades. Once provision has been made for such a dictionary structure, a sequence of gaps can be modelled as a sequence of codewords  $\{c_k\}$ , each being a reference to a dictionary entry, as represented with the *encoded stream* in the picture. Note that, for example, codeword  $c_9$  signals an exception, therefore the next symbol  $e$  is decoded using an escape mechanism.

and then apply ANS (VByte+ANS). Local variability can be instead captured by using 16 different ANS models, each selected using a 4-bit selector in the spirit of the Simple approach described in Section 3.2 (Simple+ANS). Another variant is obtained by dividing a list into blocks and encoding each block with the most suitable model, chosen among 16 possibilities according to a selected block statistic e.g., its maximum value (Packed+ANS).

### 4.3 Dictionary-based

Pibiri, Petri, and Moffat [79] show that inverted indexes can be effectively compressed using a dictionary-based approach. Their technique – named Dictionary of INTEger sequences (DINT) – builds on the observation that patterns of gaps are highly repetitive across the whole inverted index. For example, the pattern  $[1, 1, 2, 1]$  of 4 gaps can be very repetitive. Therefore, a dictionary storing the *most frequent*  $2^b$  patterns, for some  $b > 0$ , can be constructed. Note that, in general, the problem of building a dictionary that minimizes the number of output bits when sequences symbols are coded as references to its entries is NP-hard [98]. More specifically, an integer list can be modelled as a sequence of  $b$ -bit codewords, each codeword corresponding to a dictionary pattern. Fig. 6 illustrates the approach. This representation has the twofold advantage of: (1) requiring  $b$  bits to represent a pattern (thus, potentially, several integers); (2) decoding of a pattern requires just a lookup in the dictionary. In their investigation, patterns of size 1, 2, 4, 8 and 16 are considered, with  $b = 16$  to avoid bit-level manipulations and allow very fast decoding.

A detail of crucial importance is to take advantage of the presence of *runs* of 1s, hence reserving some special entries in the dictionary to encode runs of different sizes, such as 32, 64, 128, and 256. A dictionary entry must also be reserved to signal the presence of an *exception* – an integer not sufficiently frequent to be included in the dictionary and represented via an escape mechanism (e.g., Variable-Byte or a plain 32-bit integer). Moreover, compacting the dictionary has the potential of letting the dictionary fit in the processor cache, hence speeding up the decoding process thanks to a reduced number of cache misses. Lastly, once the dictionary is built, a shortest-path computation suffices to find the *optimal encoding* of a list for that specific dictionary.

Other authors have instead advocated the use of *Re-Pair* [48] to compress the gaps between the integers of inverted lists [18, 19]. Re-Pair uses a grammar-based approach to generate a dictionary of sequences of symbols. Description of this algorithm is outside the scope of this article.



## 5 FURTHER READINGS

Besides the individual papers listed in the bibliography, we mention here previous efforts in summarizing encoding techniques for integers/integer sequences. The book by Witten, Moffat, and Bell [109] is the first, to the best of our knowledge, that treats compression and indexing data as a unified problem, by presenting techniques to solve it efficiently. Fenwick [34] and Salomon [90] provide a vast and deep coverage of variable-length codes. The survey by Zobel and Moffat [113] covers more than 40 years of academic research in Information Retrieval and gives an introduction to the field, with Section 8 dealing with efficient index representations. Moffat and Turpin [68], Moffat [58], Pibiri and Venturini [82] describe several of the techniques illustrated in this article; Williams and Zobel [108], Scholer et al. [92] and Trotman [102] experimentally evaluate many of them.

Other approaches not described in this article include: an adaptation of *Front Coding* [109] for compressing text, seen as formed by quadruples holding document, paragraph, sentence, and word number [16]; the use of general-purpose compression libraries, such as ZStd<sup>3</sup> and XZ<sup>4</sup>, for encoding/decoding of inverted lists [77].

## 6 EXPERIMENTS

In this section of the article, we report on the space effectiveness and time efficiency of different inverted index representations. Specifically, space effectiveness is measured as the average number of bits dedicated to the representation of a document identifier; time efficiency is assessed in terms of the time needed to perform sequential decoding, intersection, and union of inverted lists. For the latter two operations, we focus on materializing the *full* results set, without any ranking or dynamic pruning mechanism [13, 55] being applied.

We do not aim at being exhaustive here but rather compare some selected representations and point the interested reader to the code repository at [https://github.com/jermp/2i\\_bench](https://github.com/jermp/2i_bench) for further comparisons.

**Tested index representations.** We compare the 12 different configurations, summarized in Table 9. We report some testing details of such configurations. The tested Rice implementation (Section 2.5) specifies the Rice parameter  $k$  for each block of integers, choosing the value of  $k \in [1, 4]$  giving the best space effectiveness. Two bits per block suffices to encode the value of  $k$ . Also, we write the quotient of the Rice representation of an integer in  $\gamma$  rather than in unary, as we found this to give a better space/time trade-off than regular Rice. Variable-Byte uses the SIMD-ized decoding algorithm devised by Plaisance et al. [84] and called Masked-VByte. Interpolative (BIC) uses *leftmost minimal* binary codes. The tested version of DINT uses a single packed dictionary and optimal block parsing. In Roaring, extremely dense chunks are represented with runs.

**Datasets.** We perform the experiments on the following standard test collections. Gov2 is the TREC 2004 Terabyte Track test collection, consisting in roughly 25 million .gov sites crawled in early 2004. The documents are truncated to 256 KB. ClueWeb09 is the ClueWeb 2009 TREC Category B test collection, consisting in roughly 50 million English web pages crawled between January and February 2009. CCNews is a dataset of news freely available from [CommonCrawl](https://commoncrawl.org/). Precisely, the dataset consists of the news appeared from 09/01/16 to 30/03/18.

Identifiers were assigned to documents according to the lexicographic order of their URLs [95] (see also the discussion at the beginning of Section 3). From the original collections we retain all lists whose size is larger than 4096. The postings belonging to these lists cover 93%, 94%, and 98% of the total postings of Gov2, ClueWeb09, and CCNews respectively. From the TREC 2005 and TREC

<sup>3</sup><http://www.zstd.net>

<sup>4</sup><http://tukaani.org/xz>

Table 9. The different tested index representations.

Method	Partitioned by	SIMD	Alignment	Description
VByte	cardinality	yes	byte	fixed-size partitions of 128
Opt-VByte	cardinality	yes	bit	variable-size partitions
BIC	cardinality	no	bit	fixed-size partitions of 128
$\delta$	cardinality	no	bit	fixed-size partitions of 128
Rice	cardinality	no	bit	fixed-size partitions of 128
PEF	cardinality	no	bit	variable-size partitions
DINT	cardinality	no	16-bit word	fixed-size partitions of 128
Opt-PFor	cardinality	no	32-bit word	fixed-size partitions of 128
Simple16	cardinality	no	64-bit word	fixed-size partitions of 128
QMX	cardinality	yes	128-bit word	fixed-size partitions of 128
Roaring	universe	yes	byte	single-span
Slicing	universe	yes	byte	multi-span

Table 10. The datasets used in the experiments.

(a) basic statistics				(b) TREC 2005/06 queries			
	Gov2	ClueWeb09	CCNews		Gov2	ClueWeb09	CCNews
Lists	39,177	96,722	76,474	Queries	34,327	42,613	22,769
Universe	24,622,347	50,131,015	43,530,315	2 terms	32.2%	33.6%	37.5%
Integers	5,322,883,266	14,858,833,259	19,691,599,096	3 terms	26.8%	26.5%	27.3%
Entropy of the gaps	3.02	4.46	5.44	4 terms	18.2%	17.7%	16.8%
$\lceil \log_2 \rceil$ of the gaps	1.35	2.28	2.99	5+ terms	22.8%	22.2%	18.4%

2006 Efficiency Track topics, we selected all queries whose terms are in the lexicons of the tested collection. Table 10 reports the statistics for the collections.

**Experimental setting and methodology.** Experiments are performed on a server machine equipped with Intel i9-9900K cores (@3.6 GHz), 64 GB of RAM DDR3 (@2.66 GHz) and running Linux 5 (64 bits). Each core has two private levels of cache memory: 32 KiB L1 cache (one for instructions and one for data); 256 KiB for L2 cache. A shared L3 cache spans 16,384 KiB.

The whole code is written in C++ and compiled with gcc 9.2.1 using the highest optimization setting, i.e., with compilation flags `-O3` and `-march=native`.

We build the indexes in internal memory and write the corresponding data structures to a file on disk. To perform the queries, the data structure is memory mapped from the file and a warming-up run is executed to fetch the necessary pages from disk. To test the speed of intersection and union, we use a random sampling of 1000 queries for each number of query terms from 2 to 5+ (with 5+ meaning queries with *at least* 5 terms). Each experiment was repeated 3 times to smooth fluctuations during measurements. The time reported is the average among these runs.

**Compression effectiveness.** In Table 11 we report the compression effectiveness of each method expressed as total GiB and bit-per-integer rate. The following considerations hold pretty much consistently across the three tested datasets. The most effective method is BIC with PEF being close second. Observe that both methods come very close to the entropy of gaps (with BIC being even better), as reported if Table 10a. The least effective methods are VByte and Roaring (in particular, Roaring is sensibly better than VByte on Gov2 but performs worse on the other two datasets). The representations Opt-VByte,  $\delta$ , Rice, DINT, Opt-PFor and Simple16 are all similar in space, taking

Table 11. Space effectiveness in total GiB and bits per integer, and nanoseconds per decoded integer.

Method	Gov2			ClueWeb09			CCNews		
	GiB	bits/int	ns/int	GiB	bits/int	ns/int	GiB	bits/int	ns/int
VByte	5.46	8.81	0.96	15.92	9.20	1.09	21.29	9.29	1.03
Opt-VByte	2.41	3.89	0.73	9.89	5.72	0.92	14.73	6.42	0.72
BIC	1.82	2.94	5.06	7.66	4.43	6.31	12.02	5.24	6.97
$\delta$	2.32	3.74	3.56	8.95	5.17	3.72	14.58	6.36	3.85
Rice	2.53	4.08	2.92	9.18	5.31	3.25	13.34	5.82	3.32
PEF	1.93	3.12	0.76	8.63	4.99	1.10	12.50	5.45	1.31
DINT	2.19	3.53	1.13	9.26	5.35	1.56	14.76	6.44	1.65
Opt-PFor	2.25	3.63	1.38	9.45	5.46	1.79	13.92	6.07	1.53
Simple16	2.59	4.19	1.53	10.13	5.85	1.87	14.68	6.41	1.89
QMX	3.17	5.12	0.80	12.60	7.29	0.87	16.96	7.40	0.84
Roaring	4.11	6.63	0.50	16.92	9.78	0.71	21.75	9.49	0.61
Slicing	2.67	4.31	0.53	12.21	7.06	0.68	17.83	7.78	0.69

roughly 3 – 4, 5 – 6 and 6 – 6.5 bits/int for Gov2, ClueWeb09, and CCNews respectively. The QMX and Slicing approaches stand in a middle position between the former two classes of methods.

**Sequential decoding.** Table 11 also reports the average nanoseconds spent per decoded integer, measured after decoding all lists in the index. For all the different methods, the result of decoding a list is materialized into an output buffer of 32-bit integers. Again, results are consistent across the different datasets.

The fastest methods are Roaring and Slicing thanks to their “simpler” design involving byte-aligned codes, bitmaps, and the use of SIMD instructions, allowing a value to be decoded in 0.5 – 0.7 nanoseconds. The methods Opt-VByte, QMX, and PEF are the second fastest, requiring 0.7 – 1.3 nanoseconds on average. In particular, Opt-VByte and PEF gain most of their speed thanks to the efficient decoding of dense bitmaps. The methods BIC,  $\delta$ , and Rice are the slowest as they only decode one symbol at a time (observe that BIC is almost 2 $\times$  slower than the other two because of its recursive implementation). The other mechanisms VByte, DINT, Opt-PFor and Simple16 provide similar efficiency, on average decoding an integer in 1 – 1.9 nanoseconds.

Lastly, recall that all methods – except BIC, PEF, Roaring, and Slicing – require a prefix-sum computation because they encode the gaps between the integers. In our experiments, we determined that the cost of computing the prefix-sum of the gaps is 0.5 nanoseconds per integer. This cost, sometimes, dominates that of decoding the gaps.

**Boolean AND/OR queries.** We now consider the operations of list intersection and union. Table 12 and 13 report the timings by varying the number of query terms. Fig. 7 displays the data in the tables for the ClueWeb09 dataset along space/time trade-off curves, (thus, also incorporating the space information brought by Table 11) and with the time being the “avg.” column. Almost identical shapes were obtained for the other datasets. When considering the general trade-off, especially highlighted by the plots, we see that the trend of the trade-off is the same for both intersections and unions, even across three different datasets. Therefore we can make some general points.

For methods partitioned by cardinality, the efficiency of intersection is strictly correlated to that of NextGEQ( $x$ ), an operation returning the smallest integer  $z \geq x$ ; the efficiency of union is correlated to that of sequential decoding. This is not necessarily true for Roaring and Slicing that, being partitioned by universe rather than cardinality, employ an intersection algorithm that does not use NextGEQ, nor a merging algorithm that loop through every single integer in a sequence.

Table 12. Milliseconds spent per AND query by varying the number of query terms.

Method	Gov2					ClueWeb09					CCNews				
	2	3	4	5+	avg.	2	3	4	5+	avg.	2	3	4	5+	avg.
VByte	2.2	2.8	2.7	3.3	2.8	10.2	12.1	13.7	13.9	12.5	14.0	22.4	19.7	21.9	19.5
Opt-VByte	2.8	3.1	2.8	3.2	3.0	12.2	13.3	14.0	13.6	13.3	16.0	23.2	19.6	20.3	19.8
BIC	6.8	9.7	10.4	13.2	10.0	31.7	44.2	51.5	53.8	45.3	45.6	79.7	76.9	88.8	72.8
$\delta$	4.6	6.3	6.5	8.2	6.4	20.9	28.3	33.5	34.5	29.3	28.6	50.9	48.0	55.6	45.8
Rice	4.1	5.6	5.8	7.3	5.7	19.2	25.7	30.2	31.1	26.6	26.5	46.5	43.5	50.1	41.6
PEF	2.5	3.1	2.8	3.2	2.9	12.3	13.5	14.4	13.8	13.5	17.2	24.6	21.0	21.9	21.2
DINT	2.5	3.3	3.3	4.1	3.3	11.9	14.6	16.5	17.1	15.0	16.9	27.3	24.6	28.1	24.2
Opt-PFor	2.6	3.5	3.5	4.3	3.5	12.8	15.9	18.0	18.3	16.3	16.6	27.2	24.3	27.1	23.8
Simple16	2.8	3.7	3.7	4.6	3.7	12.8	16.3	18.4	18.9	16.6	17.6	28.8	26.3	29.5	25.5
QMX	2.0	2.6	2.5	3.0	2.5	9.6	11.5	13.0	13.1	11.8	13.3	21.5	18.8	20.8	18.6
Roaring	0.3	0.5	0.7	0.8	0.6	1.5	2.5	3.1	4.3	2.9	1.1	2.0	2.6	4.1	2.5
Slicing	0.3	1.0	1.2	1.6	1.0	1.5	4.5	5.4	6.7	4.5	1.8	4.3	5.1	6.0	4.3

There is a cluster of techniques providing similar efficiency/effectiveness trade-offs, including PEF, DINT, Opt-VByte, Simple16, Opt-PFor and QMX, whereas BIC,  $\delta$  and Rice are always the slowest and dominated by the aforementioned techniques.

The procedures employed by Roaring and Slicing outperform in efficiency all techniques by a wide margin. Again, this is possible because they do not consider one-symbol-at-a-time operations, rather they rely on the intrinsic parallelism of inexpensive bitwise instructions over 64-bit words. The difference in query time between Roaring and Slicing has to be mostly attributed to the SIMD instructions that are better exploited by Roaring thanks to its simpler design. To confirm this, we performed an experiment with SIMD disabled and obtained almost identical timings to those of Slicing. However, these representations take more space than the aforementioned cluster of techniques: Slicing stands in a middle position between such cluster and Roaring.

Also observe that, for AND queries, the efficiency gap between the methods partitioned by universe and the ones partitioned by cardinality reduces when more query terms are considered. This is because the queries becomes progressively more selective (on average), hence allowing large skips to be performed by NextGEQ. On the contrary, methods partitioned by universe only skip at a coarser level, e.g., chunks containing at most  $2^{16}$  integers, therefore the cost for in-chunk calculations is always paid, even when only few integers belong to the result set. Note that this is not true for OR queries: the gap becomes progressively more evident with more query terms.

## 7 CONCLUSIONS AND FUTURE RESEARCH DIRECTIONS

The problem of introducing a compression format for sorted integer sequences, with good practical intersection/union performance, is well-studied and important, given its fundamental application to large-scale retrieval systems such as Web search engines. For that reason, inverted index compression is still a very active field of research that began several decades ago. With this article, we aimed at surveying the encoding algorithms suitable to solve the problem. However, electing a solution as the “best” one is not generally easy, rather the many space/time trade-offs available can satisfy different application requirements and the solution should always be determined by considering the actual data distribution. To this end, we also offer an experimental comparison between many of the techniques described in this article. The different space/time trade-offs assessed by this analysis are summarized by Fig. 7.

Because of the maturity reached by the state-of-the-art and the specificity of the problem, identifying future research directions is not immediate. We mention some promising ones. In

Table 13. Milliseconds spent per OR query by varying the number of query terms.

Method	Gov2					ClueWeb09					CCNews				
	2	3	4	5+	avg.	2	3	4	5+	avg.	2	3	4	5+	avg.
VByte	6.8	24.4	54.7	131.7	54.4	20.1	71.3	156.0	379.5	156.7	24.4	94.5	178.8	391.4	172.3
Opt-VByte	11.0	35.7	77.4	176.0	75.0	31.3	101.4	213.4	500.1	211.6	36.4	128.0	232.0	510.4	226.7
BIC	16.7	50.3	105.0	238.8	102.7	49.9	145.3	290.4	668.2	288.4	64.4	193.8	332.6	692.5	320.8
$\delta$	12.6	40.8	87.9	202.5	85.9	34.9	112.9	236.7	557.7	235.6	42.2	144.9	263.8	571.3	255.5
Rice	13.4	43.1	93.3	211.3	90.3	36.8	118.2	248.5	576.6	245.0	43.6	149.3	270.5	585.6	262.2
PEF	10.2	33.0	71.7	164.2	69.8	31.1	99.7	208.5	492.3	207.9	37.6	127.5	232.6	507.1	226.2
DINT	8.5	28.5	63.7	147.6	62.1	24.9	84.1	178.8	424.3	178.0	30.6	109.2	200.4	432.7	193.2
Opt-PFor	8.9	31.1	69.4	161.4	67.7	27.0	90.8	194.0	453.5	191.3	31.3	113.2	209.0	447.2	200.2
Simple16	7.8	26.2	58.3	138.2	57.6	23.7	78.0	165.5	394.7	165.5	28.7	101.5	185.3	397.8	178.4
QMX	6.6	23.8	53.4	128.1	53.0	19.7	70.0	153.2	377.9	155.2	24.0	92.6	175.2	382.4	168.6
Roaring	1.2	2.8	4.3	6.4	3.7	4.7	9.0	12.0	15.7	10.3	3.8	7.6	10.5	15.1	9.2
Slicing	1.3	4.0	6.3	9.2	5.2	5.0	12.8	18.1	25.3	15.3	5.8	12.9	17.3	23.0	14.8

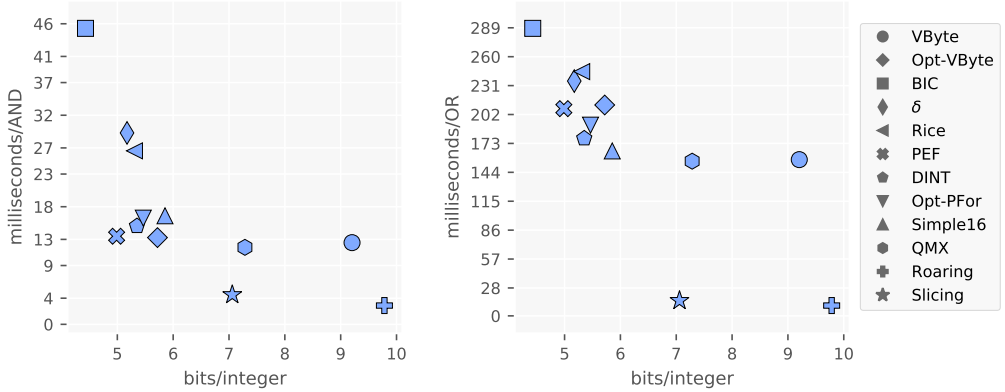


Fig. 7. Space/time trade-off curves for the ClueWeb09 dataset.

general, devising “simpler” compression formats that can be decoded with algorithms using low-latency instructions (e.g., bitwise) and with as few branches as possible, is a profitable line of research, as demonstrated by the experimentation in this article. Such algorithms favour the super-scalar execution of modern CPUs and are also suitable for SIMD instructions. Another direction could look at devising *dynamic and compressed* representations for integer sequences, able of also supporting additions and deletions. This problem is actually a specific case of the more general *dictionary problem*, which is a fundamental textbook problem. While a theoretical solution already exists with all operations supported in optimal time and compressed space [81], an implementation with good practical performance could be of great interest for dynamic inverted indexes.

## ACKNOWLEDGMENTS

The authors are grateful to Daniel Lemire, Alistair Moffat, Giuseppe Ottaviano, Matthias Petri, Sebastiano Vigna, and the anonymous referees for having carefully read earlier versions of the manuscript. Their valuable suggestions substantially improved the quality of exposition, shape, and content of the article.

This work was partially supported by the BIGDATAGRAPES (EU H2020 RIA, grant agreement N°780751), the “Algorithms, Data Structures and Combinatorics for Machine Learning” (MIUR-PRIN 2017), and the OK-INSAD (MIUR-PON 2018, grant agreement N°ARS01\_00917) projects.

## REFERENCES

- [1] Norman Abramson. 1963. *Information theory and coding*. McGraw-Hill.
- [2] Vo Ngoc Anh and Alistair Moffat. 1998. Compressed Inverted Files with Reduced Decoding Overheads. In *Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '98)*. ACM, New York, NY, USA, 290–297.
- [3] Vo Ngoc Anh and Alistair Moffat. 2005. Inverted Index Compression Using Word-Aligned Binary Codes. *Information Retrieval Journal* 8, 1 (2005), 151–166.
- [4] Vo Ngoc Anh and Alistair Moffat. 2010. Index compression using 64-bit words. *Software: Practice and Experience* 40, 2 (2010), 131–147.
- [5] Naiyong Ao, Fan Zhang, Di Wu, Douglas S Stones, Gang Wang, Xiaoguang Liu, Jing Liu, and Sheng Lin. 2011. Efficient parallel lists intersection and index compression algorithms using graphics processing units. *Proceedings of the VLDB Endowment* 4, 8 (2011), 470–481.
- [6] Alberto Apostolico and A Fraenkel. 1987. Robust transmission of unbounded strings using Fibonacci representations. *IEEE Transactions on Information Theory* 33, 2 (1987), 238–245.
- [7] Dan Blandford and Guy Blelloch. 2002. Index compression through document reordering. In *Proceedings DCC 2002. Data Compression Conference*. IEEE, 342–351.
- [8] Paolo Boldi and Sebastiano Vigna. 2004. The WebGraph framework II: Codes for the World-Wide Web. In *Data Compression Conference*. 1.
- [9] Paolo Boldi and Sebastiano Vigna. 2005. Codes for the World Wide Web. *Internet Mathematics* 2, 4 (2005), 407–429.
- [10] Abraham Bookstein and Shmuel T Klein. 1989. Construction of optimal graphs for bit-vector compression. In *Proceedings of the 13th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 327–342.
- [11] Nieves R Brisaboa, Antonio Farina, Gonzalo Navarro, and Maria F Esteller. 2003. (S, C)-dense coding: An optimized compression code for natural language text databases. In *International Symposium on String Processing and Information Retrieval*. Springer, 122–136.
- [12] Nieves R Brisaboa, Susana Ladra, and Gonzalo Navarro. 2013. DACs: Bringing direct access to variable-length codes. *Information Processing & Management* 49, 1 (2013), 392–404.
- [13] Andrei Z. Broder, David Carmel, Michael Herscovici, Aya Soffer, and Jason Y. Zien. 2003. Efficient query evaluation using a two-level retrieval process. In *Proceedings of the 12th ACM International Conference on Information and Knowledge Management*. 426–434.
- [14] Stefan Büttcher, Charles Clarke, and Gordon Cormack. 2010. *Information retrieval: implementing and evaluating search engines*. MIT Press.
- [15] Samy Chambi, Daniel Lemire, Owen Kaser, and Robert Godin. 2016. Better bitmap performance with Roaring bitmaps. *Software: practice and experience* 46, 5 (2016), 709–719.
- [16] Y. Choueka, A. S. Fraenkel, and S. T. Klein. 1988. Compression of concordances in full-text retrieval systems. In *Proceedings of the 9th annual international ACM SIGIR conference on Research and development in information retrieval*. 597–612.
- [17] David Clark. 1996. *Compact Pat Trees*. Ph.D. Dissertation. University of Waterloo.
- [18] Francisco Claude, Antonio Fariña, Miguel A. Martínez-Prieto, and Gonzalo Navarro. 2016. Universal indexes for highly repetitive document collections. *Information Systems* 61 (2016), 1–23.
- [19] F. Claude, A. Fariña, and G. Navarro. 2009. Re-Pair compression of inverted lists. *CoRR* abs/0911.3318 (2009). <http://arxiv.org/abs/0911.3318>
- [20] Intel Corporation. [last checked April 2019]. The Intel Intrinsics Guide, <https://software.intel.com/sites/landingpage/IntrinsicsGuide/>.
- [21] J Shane Culpepper and Alistair Moffat. 2005. Enhanced byte codes with restricted prefix properties. In *International Symposium on String Processing and Information Retrieval*. Springer, 1–12.
- [22] Michael Curtiss, Iain Becker, Tudor Bosman, Sergey Doroshenko, Lucian Grijincu, Tom Jackson, Sandhya Kunnatur, Soren Lassen, Philip Pronin, Sriram Sankar, Guanghao Shen, Gintaras Woss, Chao Yang, and Ning Zhang. 2013. Unicorn: A System for Searching the Social Graph. In *Proceedings of the Very Large Database Endowment*, Vol. 6. 1150–1161.
- [23] Jeffrey Dean. 2009. Challenges in building large-scale information retrieval systems: invited talk. In *Proceedings of the 2nd International Conference on Web Search and Data Mining*.



- [24] Biplob Debnath, Sudipta Sengupta, and Jin Li. 2011. Skimpystash: RAM space skimpy key-value store on flash-based storage. In *Proceedings of the 2011 ACM SIGMOD International Conference on Management of data*. ACM, 25–36.
- [25] Renaud Delbru, Stéphane Campinas, and Giovanni Tummarello. 2012. Searching web data: An entity retrieval and high-performance indexing model. *Journal of Web Semantics* 10 (2012), 33–58.
- [26] Laxman Dhulipala, Igor Kabiljo, Brian Karrer, Giuseppe Ottaviano, Sergey Pupyrev, and Alon Shalita. 2016. Compressing Graphs and Indexes with Recursive Graph Bisection. In *Proceedings of the 22nd International Conference on Knowledge Discovery and Data Mining*. 1535–1544.
- [27] Jarek Duda. 2009. Asymmetric numeral systems. CoRR abs/0902.0271 (2009). <http://arxiv.org/abs/0902.0271>
- [28] Jarek Duda. 2013. Asymmetric numeral systems: Entropy coding combining speed of Huffman coding with compression rate of arithmetic coding. CoRR abs/1311.2540 (2013). <http://arxiv.org/abs/1311.2540>
- [29] Jarek Duda, Khalid Tahboub, Neeraj J Gadgil, and Edward J Delp. 2015. The use of asymmetric numeral systems as an accurate replacement for Huffman coding. In *2015 Picture Coding Symposium (PCS)*. IEEE, 65–69.
- [30] Peter Elias. 1974. Efficient Storage and Retrieval by Content and Address of Static Files. *J. ACM* 21, 2 (1974), 246–260.
- [31] Peter Elias. 1975. Universal codeword sets and representations of the integers. *IEEE Transactions on Information Theory* 21, 2 (1975), 194–203.
- [32] Robert Mario Fano. 1949. *The transmission of information*. Massachusetts Institute of Technology, Research Laboratory of Electronics.
- [33] Robert Mario Fano. 1971. On the number of bits required to implement an associative memory. *Memorandum 61, Computer Structures Group, MIT* (1971).
- [34] Peter Fenwick. 2003. Universal codes. *Lossless Compression Handbook* (2003), 55–78.
- [35] AS Fraenkel, ST Klein, Y Choueka, and E Segal. 1986. Improved hierarchical bit-vector compression in document retrieval systems. In *Proceedings of the 9th annual international ACM SIGIR conference on Research and development in information retrieval*. 88–96.
- [36] Aviezri S Fraenkel and Shmuel T Klein. 1985. Novel compression of sparse bit-strings—preliminary report. In *Combinatorial algorithms on words*. Springer, 169–183.
- [37] Aviezri S Fraenkel and Shmuel T Klein. 1985. *Robust universal complete codes as alternatives to Huffman codes*. Department of Applied Mathematics, Weizmann Institute of Science.
- [38] Robert Gallager and David Van Voorhis. 1975. Optimal source codes for geometrically distributed integer alphabets (corresp.). *IEEE Transactions on Information theory* 21, 2 (1975), 228–230.
- [39] Jonathan Goldstein, Raghu Ramakrishnan, and Uri Shaft. 1998. Compressing Relations and Indexes. In *Proceedings of the 14th International Conference on Data Engineering*. 370–379.
- [40] Solomon Golomb. 1966. Run-length encodings. *IEEE Transactions on Information Theory* 12, 3 (1966), 399–401.
- [41] Rodrigo González, Szymon Grabowski, Veli Mäkinen, and Gonzalo Navarro. 2005. Practical implementation of rank and select queries. In *Workshop on Efficient and Experimental Algorithms*. 27–38.
- [42] Vagelis Hristidis, Yannis Papakonstantinou, and Luis Gravano. 2003. Efficient IR-Style Keyword Search over Relational Databases. In *Proceedings 2003 VLDB Conference*. Elsevier, 850–861.
- [43] David A Huffman. 1952. A method for the construction of minimum-redundancy codes. *Proceedings of the IRE* 40, 9 (1952), 1098–1101.
- [44] Guy Jacobson. 1989. *Succinct Static Data Structures*. Ph.D. Dissertation. Carnegie Mellon University.
- [45] Matti Jakobsson. 1978. Huffman Coding in Bit-Vector Compression. *Inf. Process. Lett.* 7, 6 (1978), 304–307.
- [46] Aaron Kiely. 2004. Selecting the Golomb parameter in Rice coding. *IPN progress report* 42 (2004), 159.
- [47] Leon Gordon Kraft. 1949. *A device for quantizing, grouping, and coding amplitude-modulated pulses*. Ph.D. Dissertation. Massachusetts Institute of Technology.
- [48] N. Jesper Larsson and Alistair Moffat. 1999. Offline Dictionary-Based Compression. In *Data Compression Conference*. 296–305.
- [49] Debra A. Lelewer and Daniel S. Hirschberg. 1987. Data Compression. *ACM Comput. Surv.* 19, 3 (Sept. 1987), 261–296.
- [50] Daniel Lemire and Leonid Boytsov. 2015. Decoding billions of integers per second through vectorization. 45, 1 (2015), 1–29.
- [51] Daniel Lemire, Owen Kaser, Nathan Kurz, Luca Deri, Chris O’Hara, François Saint-Jacques, and Gregory Ssi-Yan-Kai. 2018. Roaring bitmaps: Implementation of an optimized software library. *Software: Practice and Experience* 48, 4 (2018), 867–895.
- [52] Daniel Lemire, Nathan Kurz, and Christoph Rupp. 2018. Stream-VByte: faster byte-oriented integer compression. *Inform. Process. Lett.* 130 (2018), 1–6.
- [53] Daniel Lemire, Gregory Ssi-Yan-Kai, and Owen Kaser. 2016. Consistently faster and smaller compressed bitmaps with roaring. *Software: Practice and Experience* 46, 11 (2016), 1547–1569.
- [54] Veli Mäkinen and Gonzalo Navarro. 2007. Rank and select revisited and extended. *Theoretical Computer Science* 387, 3 (2007), 332–347.

- [55] Antonio Mallia, Giuseppe Ottaviano, Elia Porciani, Nicola Tonellotto, and Rossano Venturini. 2017. Faster BlockMax WAND with Variable-sized Blocks. In *Proceedings of the International ACM Conference on Research and Development in Information Retrieval*. 625–634.
- [56] Christopher Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. *Introduction to Information Retrieval*. Cambridge University Press.
- [57] Brockway McMillan. 1956. Two inequalities implied by unique decipherability. *IRE Transactions on Information Theory* 2, 4 (1956), 115–116.
- [58] Alistair Moffat. 2016. Compressing Integer Sequences. In *Encyclopedia of Algorithms*. 407–412.
- [59] Alistair Moffat. 2019. Huffman Coding. *ACM Comput. Surv.* 52, 4, Article 85 (2019), 35 pages.
- [60] Alistair Moffat and Vo Ngoc Anh. 2005. Binary codes for non-uniform sources. In *Data Compression Conference*. IEEE, 133–142.
- [61] Alistair Moffat and Vo Ngoc Anh. 2006. Binary codes for locally homogeneous sequences. *Inform. Process. Lett.* 99, 5 (2006), 175–180.
- [62] Alistair Moffat, Radford M. Neal, and Ian H. Witten. 1998. Arithmetic Coding Revisited. *ACM Trans. Inf. Syst.* 16, 3 (July 1998), 256–294.
- [63] Alistair Moffat and Matthias Petri. 2017. ANS-Based Index Compression. In *Proceedings of the ACM on Conference on Information and Knowledge Management*. 677–686.
- [64] Alistair Moffat and Andrew Turpin. 2018. Index Compression Using Byte-Aligned ANS Coding and Two-Dimensional Contexts. In *Proceedings of the 11th ACM International Conference on Web Search and Data Mining*. 405–413.
- [65] Alistair Moffat and Lang Stuiver. 1996. Exploiting Clustering in Inverted File Compression. In *Data Compression Conference*. 82–91.
- [66] Alistair Moffat and Lang Stuiver. 2000. Binary Interpolative Coding for Effective Index Compression. *Information Retrieval Journal* 3, 1 (2000), 25–47.
- [67] Alistair Moffat and Andrew Turpin. 1997. On the implementation of minimum redundancy prefix codes. *IEEE Transactions on Communications* 45, 10 (1997), 1200–1207.
- [68] Alistair Moffat and Andrew Turpin. 2002. *Compression and coding algorithms*. Springer Science & Business Media.
- [69] Alistair Moffat and Justin Zobel. 1992. Parameterised compression for sparse bitmaps. In *Proceedings of the 15th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 274–285.
- [70] Alistair Moffat and Justin Zobel. 1996. Self-indexing inverted files for fast text retrieval. *ACM Transactions on Information Systems (TOIS)* 14, 4 (1996), 349–379.
- [71] Leonardo of Pisa (known as Fibonacci). 1202. *Liber Abaci*.
- [72] Giuseppe Ottaviano, Nicola Tonellotto, and Rossano Venturini. 2015. Optimal Space-time Tradeoffs for Inverted Indexes. In *International ACM Conference on Web Search and Data Mining*. 47–56.
- [73] Giuseppe Ottaviano and Rossano Venturini. 2014. Partitioned Elias-Fano Indexes. In *Proceedings of the 37th International Conference on Research and Development in Information Retrieval*. 273–282.
- [74] Rasmus Pagh. 2001. Low redundancy in static dictionaries with constant query time. *SIAM J. Comput.* 31, 2 (2001), 353–363.
- [75] Athanasios Papoulis. 1991. *Probability, Random Variables, and Stochastic Processes* (3rd ed.). McGraw-Hill.
- [76] Richard Clark Pasco. 1976. *Source coding algorithms for fast data compression*. Ph.D. Dissertation. Stanford University Palo Alto, CA.
- [77] Matthias Petri and Alistair Moffat. 2018. Compact inverted index storage using general-purpose compression libraries. *Software: Practice and Experience* 48, 4 (2018), 974–982.
- [78] Giulio Ermanno Pibiri. 2019. On Slicing Sorted Integer Sequences. *CoRR abs/1907.01032* (2019). <http://arxiv.org/abs/1907.01032>
- [79] Giulio Ermanno Pibiri, Matthias Petri, and Alistair Moffat. 2019. Fast Dictionary-Based Compression for Inverted Indexes. In *International ACM Conference on Web Search and Data Mining*. 9.
- [80] Giulio Ermanno Pibiri and Rossano Venturini. 2017. Clustered Elias-Fano indexes. *ACM Transactions on Information Systems* 36, 1, Article 2 (2017), 33 pages.
- [81] Giulio Ermanno Pibiri and Rossano Venturini. 2017. Dynamic Elias-Fano Representation. In *Proceedings of the 28-th Annual Symposium on Combinatorial Pattern Matching*. 30:1–30:14.
- [82] Giulio Ermanno Pibiri and Rossano Venturini. 2018. Inverted Index Compression. *Encyclopedia of Big Data Technologies* (2018), 1–8.
- [83] Giulio Ermanno Pibiri and Rossano Venturini. 2019. On Optimally Partitioning Variable-Byte Codes. *IEEE Transactions on Knowledge and Data Engineering* (2019), 1–12.
- [84] Jeff Plaisance, Nathan Kurz, and Daniel Lemire. 2015. Vectorized VByte Decoding. In *International Symposium on Web Algorithms*.

- [85] Vijayshankar Raman, Lin Qiao, Wei Han, Inderpal Narang, Ying-Lin Chen, Kou-Horng Yang, and Fen-Ling Ling. 2007. Lazy, adaptive rid-list intersection, and its application to index anding. In *Proceedings of the 2007 ACM SIGMOD international conference on Management of data*. ACM, 773–784.
- [86] Robert Rice. 1991. Some practical universal noiseless coding techniques, part 3, module PSL14, K+. *Jet Propulsion Laboratory, JPL Publication* 91, 3 (1991), 132.
- [87] Robert Rice and J. Plaunt. 1971. Adaptive Variable-Length Coding for Efficient Compression of Spacecraft Television Data. *IEEE Transactions on Communications* 16, 9 (1971), 889–897.
- [88] JJ Rissanen. 1979. Arithmetic codings as number representations. *Acta Polytech. Scand. Math* 31 (1979), 44–51.
- [89] Jorma J Rissanen. 1976. Generalized Kraft inequality and arithmetic coding. *IBM Journal of research and development* 20, 3 (1976), 198–203.
- [90] David Salomon. 2007. *Variable-length Codes for Data Compression*. Springer.
- [91] Benjamin Schlegel, Rainer Gemulla, and Wolfgang Lehner. 2010. Fast integer compression using SIMD instructions. In *Proceedings of the Sixth International Workshop on Data Management on New Hardware*. ACM, 34–40.
- [92] Falk Scholer, Hugh E Williams, John Yiannis, and Justin Zobel. 2002. Compression of inverted indexes for fast query evaluation. In *Proceedings of the 25th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 222–229.
- [93] Claude Elwood Shannon. 1948. A mathematical theory of communication. *Bell system technical journal* 27, 3 (1948), 379–423.
- [94] Wann-Yun Shieh, Tien-Fu Chen, Jean Jyh-Jiun Shann, and Chung-Ping Chung. 2003. Inverted file compression through document identifier reassignment. *Information processing & management* 39, 1 (2003), 117–131.
- [95] Fabrizio Silvestri. 2007. Sorting Out the Document Identifier Assignment Problem. In *Proceedings of the 29th European Conference on IR Research*. 101–112.
- [96] Fabrizio Silvestri and Rossano Venturini. 2010. VSEncoding: Efficient Coding and Fast Decoding of Integer Lists via Dynamic Programming. In *Proceedings of the 19th International Conference on Information and Knowledge Management*. 1219–1228.
- [97] Alexander Stepanov, Anil Gangolli, Daniel Rose, Ryan Ernst, and Paramjit Oberoi. 2011. SIMD-based decoding of posting lists. In *Proceedings of the 20th International Conference on Information and Knowledge Management*. 317–326.
- [98] J. A. Storer and T. G. Szymanski. 1982. Data compression via textual substitution. 29, 4 (1982), 928–951.
- [99] Jukka Teuhola. 1978. A compression method for clustered bit-vectors. *Information processing letters* 7, 6 (1978), 308–311.
- [100] Jukka Teuhola. 2008. Tournament coding of integer sequences. *Comput. J.* 52, 3 (2008), 368–377.
- [101] Larry H Thiel and HS Heaps. 1972. Program design for retrospective searches on large data bases. *Information Storage and Retrieval* 8, 1 (1972), 1–20.
- [102] Andrew Trotman. 2003. Compressing inverted files. *Information Retrieval* 6, 1 (2003), 5–19.
- [103] Andrew Trotman. 2014. Compression, SIMD, and postings lists. In *Proceedings of the 2014 Australasian Document Computing Symposium*. ACM, 50.
- [104] Andrew Trotman and Kat Lilly. 2018. Elias Revisited: Group Elias SIMD Coding. In *Proceedings of the 23rd Australasian Document Computing Symposium*. ACM, 4.
- [105] Peter van Emde Boas. 1975. Preserving Order in a Forest in less than Logarithmic Time. In *Proceedings of the 16-th Annual Symposium on Foundations of Computer Science*. 75–84.
- [106] Peter van Emde Boas. 1977. Preserving Order in a Forest in Less Than Logarithmic Time and Linear Space. *Inform. Process. Lett.* 6, 3 (1977), 80–82.
- [107] Sebastiano Vigna. 2013. Quasi-succinct indices. In *Proceedings of the 6th ACM International Conference on Web Search and Data Mining*. 83–92.
- [108] Hugh E Williams and Justin Zobel. 1999. Compressing integers for fast file access. *Comput. J.* 42, 3 (1999), 193–201.
- [109] Ian Witten, Alistair Moffat, and Timothy Bell. 1999. *Managing gigabytes: compressing and indexing documents and images* (2nd ed.). Morgan Kaufmann.
- [110] Ian H Witten, Radford M Neal, and John G Cleary. 1987. Arithmetic coding for data compression. *Commun. ACM* 30, 6 (1987), 520–540.
- [111] Hao Yan, Shuai Ding, and Torsten Suel. 2009. Inverted index compression and query processing with optimized document ordering. In *Proceedings of the 18th International Conference on World Wide Web*. 401–410.
- [112] J. Zhang, X. Long, and T. Suel. 2008. Performance of compressed inverted list caching in search engines. In *International World Wide Web Conference (WWW)*. 387–396.
- [113] Justin Zobel and Alistair Moffat. 2006. Inverted files for text search engines. *Comput. Surveys* 38, 2 (2006), 1–56.
- [114] Marcin Zukowski, Sándor Héman, Niels Nes, and Peter Boncz. 2006. Super-Scalar RAM-CPU Cache Compression. In *Proceedings of the 22nd International Conference on Data Engineering*. 59–70.