# The Efficiency of the ANS Entropy Encoding

#### DMITRY KOSOLOBOV

Ural Federal University, Ekaterinburg, Russia dkosolobov@mail.ru

#### Abstract

The Asymmetric Numeral Systems (ANS) encoder invented by Duda in 2009 is an entropy encoder that had an immense impact on the data compression, substituting arithmetic and Huffman coding. The optimality of the ANS was studied by Duda and others but the precise asymptotic behaviour of its redundancy (in comparison to the source entropy) was not completely understood. In this paper we establish an optimal bound on the redundancy for the tabled ANS (tANS), which is the most popular ANS variant used in practice. Let  $a_1, a_2, \ldots, a_n$  be a sequence of letters from an alphabet  $\{0, 1, \ldots, \sigma-1\}$  such that each letter a occurs in it  $f_a$  times and  $n=2^r$ , for an integer r. The tANS encoder using Duda's "precise initialization" to fill the tANS table transforms this sequence into a bit string of the following length (the table of frequencies is not included in the encoding size):

$$\sum_{a \in [0..\sigma)} f_a \cdot \log \frac{n}{f_a} + O(\sigma + r),$$

where the  $O(\sigma+r)$  term can be upper bounded by  $\sigma\log e+r$ . The r-bit additive term is an artifact of the initial encoder state, in dispensable to ANS; the rest incurs a redundancy of  $O(\frac{\sigma}{n})$  bits per letter, which improves the previously known bound  $O(\frac{\sigma}{n}\log n)$ . We complement this upper bound with a series of examples showing that an  $\Omega(\sigma+r)$  redundancy is indeed necessary when  $\sigma>n/3$ ; the  $\Omega(\sigma+r)$  in our construction is greater than  $\frac{\sigma-1}{4}+r-2$ . We argue that similar examples can be described for any reasonable methods that distribute letters in the tANS table using only the knowledge about frequencies. Thus, we disprove Duda's conjecture that the redundancy is  $O(\frac{\sigma}{n^2})$  bits per letter.

We also propose a new variant of the range ANS (rANS), called rANS with fixed accuracy, that is parameterized by an integer  $k \geq 1$ . In this variant the division operation, which is unavoidable in rANS encoders, is performed only in cases when its result belongs to the range  $[2^k...2^{k+1})$ . Therefore, the division can be computed by faster methods provided k is small. We upperbound the redundancy for the rANS with fixed accuracy k by  $\frac{n}{2^k-1}\log e+r$ .

The unconventional way by which we introduce the ANS encoders might be interesting by itself. **Keywords:** asymmetric numeral systems, ANS, tANS, rANS, finite state entropy, FSE, arithmetic coding, entropy, range coding, redundancy

# 1 Introduction

Asymmetric numeral systems (ANS) is a class of entropy encoders invented by Duda in 2009 [11, 12, 14]. These encoders had a huge impact on the data compression by providing the same rates of compression as the arithmetic coding [18, 19, 21] while being as fast as the Huffman coding [17] (and even faster in some scenarios). Since the invention of ANS and the emergence of its efficient implementation by Collet [7], several high performance compressors based on ANS appeared [1, 6, 16] and it was integrated in some modern media formats [2, 3]. The theoretical community also contributed to the study of ANS in a series of works [4, 13, 20, 23, 24, 25], though less actively than the practitioners [5, 8, 15, 16].

The primary focus of the theoretical analysis for entropy encoders is in the estimation of the redundancy, the difference between the number of bits produced by the encoder and the information theoretic entropy lower bound. The results of the present paper are twofold: first, we establish tight asymptotic upper and

lower bounds on the redundancy of the most popular variant of ANS, called the tabled ANS (tANS) or, sometimes, the Finite State Entropy (FSE); second, we introduce and analyze a novel variant of the range ANS (rANS), another version of ANS used in practice.

An entropy encoder receives as its input a sequence of n numbers from a set  $\{0, 1, \ldots, \sigma-1\}$  and transforms it into a bit string. Duda analyzed his tANS encoder and found that its redundancy is  $O(\frac{\sigma}{n}\log n)$  bits per letter<sup>1</sup> provided an appropriate initialization is used for tANS tables (details are explained in the sequel); see Equation (40) in [11] (his estimation, however, is somewhat imprecise and in a slightly different setting). Based on experimental evaluations, Duda conjectured that the tight bound for the redundancy is  $O(\frac{\sigma}{n^2})$  bits per letter [12]. Yokoo and Dubé [25] investigated the same problem in more rigorous terms and closer to our setting and showed that the redundancy per letter vanishes as the length n tends to infinity while  $\sigma$  is fixed (however, they have some questionable assumptions in their derivations).

In this paper we prove that the redundancy for tANS is  $O(\frac{\sigma}{n})$  bits per letter. We complement this upper bound by a series of examples showing that it is asymptotically tight when  $\sigma > n/3$ . As in the works cited above, here we did not include in this bound the r redundant bits that are always produced by ANS encoders in the worst case (it is an artifact of the initial state of the encoder). After uncovering the constant under the big-O and including this r-bit term, the upper bound for the tANS redundancy that we establish can be expressed as  $\sigma \log e + r$  bits over all letters (not per letter). Our lower bound examples show that a redundancy of  $\frac{\sigma-1}{4} + r - 2$  bits is attainable, for  $\sigma \approx n/3$ . An important part of tANS is the initialization of its internal tables, which has a significant impact on the compression performance. The examples, however, work for a wide class of initialization algorithms so that the same redundancy can be observed for any adequate algorithm that generates the tANS tables using only the known frequencies of letters without processing the sequence itself (however, methods that analyze the sequence might avoid this effect [10, 13]).

The second contribution of the present paper is a modification of the range ANS (rANS). The rANS is another variant of ANS invented by Duda, which is noticeably slower in practice than the tANS. However, rANS has a number of advantages in some use scenarios because of which it was favoured by some practitioners (and because it is easier to learn). The main advantage of rANS is that it does not need tANS tables. Due to the less wasteful use of memory, the rANS is more suitable for (pseudo)adaptive compression or when several streams of data are encoded simultaneously [5, 15, 16]. The inherent problem of rANS that slows down it significantly is the necessity to perform the integer division during the encoding, which is expensive on modern processors. The proposed modification of rANS, which we call rANS with fixed accuracy, tries to mitigate this issue making the encoder faster while preserving the same good properties of rANS.

The new rANS takes as a regulated parameter an integer  $k \ge 1$ . It is guaranteed that the division can be performed only in cases when its result belongs to the range  $[2^k..2^{k+1})$ . Thus, the division can be computed by faster methods provided k is small. We upperbound the redundancy for the rANS with fixed accuracy k by  $\frac{n}{2^k-1}\log e+r$ . The new rANS variant is faster than the standard rANS but is still not as fast as tANS: in our experiments, the rANS with fixed accuracy k=3 was 1.8 slower than tANS implemented by Collet [7] (however, we did not use the technique of interleaving streams by Giesen [16] in our implementation). We believe that the rANS with fixed accuracy might be more suitable for hardware implementations, where its restricted division operation can be sped up more efficiently using versatile capabilities for parallelization.

The paper is organized as follows. After short preliminaries, Section 3 introduces a simple variant of the range ANS and, based on it, the tabled ANS. Section 4 presents a tight analysis of the redundancy for the tabled ANS, upper and lower bounds. Section 5 describes a novel variant of range ANS and analyzes it. We conclude with some open problems in Section 6.

### 2 Preliminaries

A letter is an element of a finite set of integers called an alphabet. We consider sequences of letters. The sequences are also sometimes called strings. Integer intervals are denoted by  $[i..j] = \{k \in \mathbb{Z} : i \leq k \leq j\}$ 

<sup>&</sup>lt;sup>1</sup>All logarithms in the paper are in base two.

and  $[i..j] = [i..j] \setminus \{j\}$ . Fix an alphabet  $[0..\sigma)$  with  $\sigma$  letters and a sequence of letters  $a_1, a_2, \ldots, a_n$  from it. For each letter a, denote by  $f_a$  the number of occurrences of a in the sequence. The number  $f_a$  is called the frequency of a. The number  $f_a/n$  is the empirical probability of a. An entropy encoder transforms this sequence into a sequence of bits that is then transmitted to a decoder. The encoder usually also transmits to the decoder a table of frequencies. However, the problem of the storage for the table is out of the scope of the present paper and we focus only on the encoding for the sequence itself implicitly assuming that both sides know the length n of the sequence and the frequencies of letters. It is well known that any encoder in the worst case should spend at least the following number of bits for the sequence encoding (even under the assumption that the table of frequencies is known), which the entropy information theoretic lower bound [9]:

$$\sum_{a \in [0..\sigma)} f_a \log \frac{n}{f_a},\tag{1}$$

assuming that  $f_a \log \frac{n}{f_a} = 0$  whenever  $f_a = 0$ . We call (1) the *entropy formula* (though the empirically calculated entropy itself is this value divided by n). The difference between the number of bits produced by an encoder and the optimal number of bits from (1) is called a *redundancy*. The redundancy is the primary focus of our analysis of ANS.

For each letter a, denote  $c_a = \sum_{a' \in [0..a)} f_{a'}$ . Typically, encoders stores two arrays of  $\sigma$  numbers  $c_a$  and  $f_a$  and some additional tables necessary for encoding. The ANS encoders are not an exception.

Our analysis utilizes the following well-known logarithmic inequalities:  $\log(1+x) \le x \log e$  and  $\log(1-x) \ge \frac{-x}{1-x} \log e = \frac{-\log e}{1/x-1}$ .

#### 3 ANS Encoders

On a high level, an ANS encoder can be described as a data structure that maintains a positive integer w that can be modified by two stack operations " $w = \operatorname{push}(w, a)$ " and " $(w, a) = \operatorname{pop}(w)$ ": push encodes a letter a into the number w and returns the modified value for w, and pop performs a reverse operation decoding the letter a and restoring the old value for w. Given a sequence of letters  $a_1, a_2, \ldots, a_n$ , the encoder pushes them consecutively and transmits the resulting integer w to a decoder, which in turn retrieves the sequence from w performing pop operations. Note that the letters appear in the reverse order during the decoding, which is a distinctive feature of ANS. For their correct coordinated work, both encoder and decoder should receive in advance the same table of (approximate) frequencies of letters in the sequence  $a_1, a_2, \ldots, a_n$ .

For didactic reasons, ANS encoders are usually first described in their unbounded form in which they operate on very long integers w, which is unrealistic. In practice, streaming ANS encoders are used, which maintain the value of w during the construction within a fixed range fitting into a machine word and they store all excessive bits in an output buffer by performing a "renormalization" (akin to arithmetic encoders). However, in what follows, we significantly diverge from this standard way of explanation for ANS and introduce only the streaming ANS but in an unconventional manner; the unbounded ANS variant is not discussed at all.

#### 3.1 Basic Range ANS

Suppose that we are to encode a sequence of letters  $a_1, a_2, \ldots, a_n$  from the alphabet  $[0..\sigma)$ . For each letter  $a \in [0..\sigma)$ , denote by  $f_a$  its frequency in the sequence (i.e., its empirical probability is  $\frac{f_a}{n}$ ). We assume that n is a power of two, i.e.,  $n=2^r$  for an integer  $r \geq 0$ . It is a standard assumption for both ANS and arithmetic coding that simplifies implementations. If the length of the sequence is not a power of two, then either the real probabilities of letters are approximated with numbers  $f_a/2^r$  or the sequence is split into chunks whose lengths are powers of two; we discuss the former case briefly in the end of Section 4.1 but most details are omitted as they are not in the scope of this paper.

The streaming ANS encoder reads the letters  $a_1, a_2, \ldots, a_n$  from left to right and, after processing  $a_1, a_2, \ldots, a_i$ , encodes the processed part into a positive integer  $w_i$ . Initially,  $w_0 = 2^r$ ; the choice for  $w_0$  is somewhat arbitrary, the only necessary condition is  $w_0 \geq 2^r$ . To encode a new letter  $a_{i+1}$ , we

transform the number  $w_i$  into  $w_{i+1}$  by increasing the value stored in the r+1 highest bits of  $w_i$ . Thus,  $w_0 < w_1 < \cdots < w_{i+1}$ . According to our terminology, we have  $w_{i+1} = \mathsf{push}(w_i, a_{i+1})$ . It is instructive to image the number  $w_i$  as a bit stream written from the highest bit (which is always 1) to lowest bits. We replace at most r+1 highest bits with a new larger value; see Figure 1.

```
w_i = \underbrace{\frac{10010}{011001000010110}}_{t+1 \text{ bits}} 011001000010110
```

Figure 1: A transformation of the number  $w_i$  into  $w_{i+1}$ .

The numbers  $w_i$  might be very long. However, since only the highest r+1 bits of  $w_i$  matter for the encoder, all lower bits can be stored in an output buffer. The integer r is chosen so that the r+1 bits can be stored into one machine word. For simplicity of the exposition, we omit this technical detail and continue to discuss the numbers  $w_i$  as if they were stored explicitly.

The goal is to encode the whole sequence optimally so that, ideally, the final number  $w_n$  occupies  $\sum_a f_a \log \frac{n}{f_a}$  bits. Intuitively, one can achieve this by encoding each letter a into  $\log \frac{n}{f_a} = r - \log f_a$  bits. During the processing of  $a = a_{i+1}$ , we could have achieved this by replacing the highest  $\log f_a + 1$  bits of  $w_i$  with r+1 bits that store a new larger value; then, the number of bits in  $w_{i+1}$  and  $w_i$  will differ by  $r+1-(\log f_a+1)=\log \frac{n}{f_a}$ . But the number  $\log f_a$  is not integer in general. Therefore, instead, we replace either  $\lfloor \log f_a \rfloor + 1$  or  $\lfloor \log f_a \rfloor + 2$  highest bits of  $w_i$  with new r+1 bits. As a result, the number of bits in  $w_{i+1}$  increases by either  $r-\lfloor \log f_a \rfloor$  or  $r-\lfloor \log f_a \rfloor -1$ , which is approximately  $\log \frac{n}{f_a}$ . The cumulative growth of  $w_n$  may approach the optimal  $\log \frac{n}{f_a}$  bits per letter a on average if the case when  $\lfloor \log f_a \rfloor + 2$  bits are replaced happens more often. (Note that when  $f_a \geq n/2 = 2^{r-1}$ , we have  $\lfloor \log f_a \rfloor = r-1$  and the number of bits in  $w_i$  sometimes might not change at all; but the content will change.)

What is the content of the r+1 new highest bits in  $w_{i+1}$  and how do we decide whether  $\lfloor \log f_a \rfloor + 1$  or  $\lfloor \log f_a \rfloor + 2$  highest bits of  $w_i$  will be replaced?

Denote by x' the value stored in the highest r+1 bits of  $w_{i+1}$ , i.e.,  $w_{i+1} = x' \cdot 2^{\lfloor \log w_{i+1} \rfloor - r} + \Delta$ , where  $0 \le \Delta < 2^{\lfloor \log w_{i+1} \rfloor - r}$  and  $2^r \le x' < 2^{r+1}$ . Denote by x the value of the highest bits of  $w_i$  that were replaced with x', i.e.,  $w_i = x \cdot 2^{\lfloor \log w_{i+1} \rfloor - r} + \Delta$ , where  $x \ge 1$ . (In the example of Figure 1 x and x' are emphasized and  $\Delta$  is a common part of  $w_i$  and  $w_{i+1}$ .) The scheme must be reversible: the number x' must provide sufficient information for the decoder to restore the letter  $a_{i+1}$  and the number x. Since  $\sigma \le \sum_a f_a = 2^r$ , there is enough place to encode  $a_{i+1}$  into the r lowest bits of x' (the highest, (r+1)th, bit of x' is 1). "Excessive"  $2^r - \sigma$  possible values of x' will be used to restore the number x, the replaced highest bits of  $w_i$ .

The r lowest bits of x' can store any number from the range  $[0..2^r)$ . The encoder chooses one of these numbers depending on the values of x and  $a_{i+1}$ . We distribute the numbers from  $[0..2^r)$  among the letters according to their frequencies: for each letter a, denote  $c_a = \sum_{a' \in [0..a)} f_{a'}$ ; the subrange of values  $[c_a..c_{a+1})$  is allocated for a. Hence, given a letter  $a = a_{i+1}$ , we have a room in the range that can store any number from  $[0..f_a)$  and this information should suffice to restore the replaced number x from x'. As an example, Figure 2 depicts a distribution of the range  $2^r = 2^4$  among letters a, b, c, d with frequencies 3, 5, 6, 2, respectively.

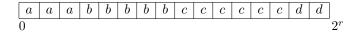


Figure 2: A distribution of letters in a range of length  $2^r$ .

Let  $a=a_{i+1}$ . Denote by  $x_1$  and  $x_2$  the values stored in, respectively, the highest  $\lfloor \log f_a \rfloor + 1$  and the highest  $\lfloor \log f_a \rfloor + 2$  bits of  $w_i$ . We assume that  $x=x_1$  if  $x_1 \geq f_a$ , and  $x=x_2$  otherwise. Thus, the condition  $x_1 \geq f_a$  determines whether  $\lfloor \log f_a \rfloor + 1$  or  $\lfloor \log f_a \rfloor + 2$  highest bits of  $w_i$  will be used for x. Note that  $2^{\lfloor \log f_a \rfloor} \leq x_1 < 2 \cdot f_a$  and  $x_2 \leq 2 \cdot x_1 + 1$ . Therefore, we have  $f_a \leq x < 2 \cdot f_a$  and, thus,  $x \mod f_a = x - f_a$ .

The latter can be clearly seen in Figure 3: the bits occupied by  $x_1$  and  $x_2$  are emphasized, respectively, on the left and right; note that the number  $f_a$  occupies  $\lfloor \log f_a \rfloor + 1$  bits.

$$w_i = \frac{10110}{1001000100101000010110}$$
  $w_i = \frac{100100}{1001000010110}$   $f_a = 10011$   $f_a = 10011$   $x \mod f_a = 11$   $x \mod f_a = 10001$ 

Figure 3: Two cases: x occupies  $\lfloor \log f_a \rfloor + 1$  (left) or  $\lfloor \log f_a \rfloor + 2$  (right) highest bits of  $w_i$ .

In order to restore the value x from x', it suffices to encode somehow the value x mod  $f_a$  into x': then, x is restored as  $x = f_a + (x \mod f_a)$ . To this end, the range  $[c_a...c_{a+1})$  allocated for the letter a has exactly enough room. Thus, we have the following transformation to encode a:

$$x' = 2^r + c_a + (x \bmod f_a), \text{ where } f_a \le x < 2 \cdot f_a.$$

With this approach, the decoder should perform a reverse transformation for the r+1 highest bits of the number  $w_{i+1}$  in order to restore  $w_i$ . The decoding is straightforward:

$$x = f_a + (x' \mod 2^r) - c_a$$
, where  $c_a \le x' \mod 2^r < c_{a+1}$ . (3)

The decoded letter a is determined by examining to which range  $[c_a..c_{a+1})$  the number  $x' \mod 2^r$  belongs. This is how the operation " $(w_i, a_{i+1}) = \mathsf{pop}(w_{i+1})$ " is performed. It remains to observe that x' > x since  $x' \ge 2^r + (x \mod f_a) > f_a + (x \mod f_a) = x$ . Therefore, the number  $w_{i+1}$  indeed is larger than  $w_i$ .

The described scheme is the simplest form of the so-called range ANS (rANS). As will be seen later, the size of  $w_n$  in bits can be bounded by  $\sum_a f_a \log \frac{n}{f_a} + O(n)$ . The redundancy O(n) is quite significant for many applications and, indeed, the described scheme does not perform well in practice. The following section describes an additional "shuffling" step in the encoder that fixes this issue. There exists, however, a more elaborate version of the range ANS that does not have such a problem and works well without shuffling; we discuss it briefly in Section 5 where we also present a novel variant of the range ANS with good compression guarantees.

#### 3.2 Shuffling and Tabled ANS

The idea of the shuffling step enhancing the simple range ANS is to shuffle the lower bits of x' in a random-like fashion so that the scheme (2) is changed as follows:

$$x' = 2^r + \mathsf{shuffle}[c_a + (x \bmod f_a)], \text{ where } f_a \le x < 2 \cdot f_a. \tag{4}$$

The array shuffle[0..2<sup>r</sup>-1] is a permutation of the range [0..2<sup>r</sup>) but it is not entirely random: in order to guarantee the inequality x' > x (which implies  $w_{i+1} > w_i$ ), it must satisfy the following property:

$$\mathsf{shuffle}[c_a + i] < \mathsf{shuffle}[c_a + j], \text{ whenever } 0 \le i < j < f_a. \tag{5}$$

Due to this condition, we have x' > x since  $x' \ge 2^r + (x \mod f_a) > f_a + (x \mod f_a) = x$ . It is convenient to view shuffle as defined via an additional array range $[0..2^r-1]$  that stores an (arbitrary) permutation of the array of letters from Figure 2 in which every letter a occurs exactly  $f_a$  times; then, for  $i \in [0..f_a)$ , shuffle $[c_a + i]$  is equal to the index of the (i + 1)th occurrence of the letter a in range (see Figure 4). Thus, to define shuffle, it suffices to initialize the array range; we will implicitly imply this relation in the sequel when the initialization of shuffle is discussed.

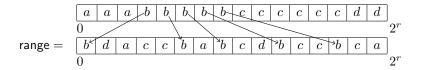


Figure 4: A shuffled distribution of letters. Here, we have  $\mathsf{shuffle}[c_b] = 0$ ,  $\mathsf{shuffle}[c_b + 1] = 5$ ,  $\mathsf{shuffle}[c_b + 2] = 7$ .

The decoding procedure (3) performs a reverse transformation in an obvious way:

$$x = f_a + \mathsf{unshuffle}[x' \bmod 2^r] - c_a$$
, where  $\mathsf{range}[x' \bmod 2^r] = a$ .

Note the use of the array range in the decoder to determine the letter a. The array unshuffle is the inverse of shuffle such that unshuffle[shuffle[z]] = z, for any  $z \in [0..2^r)$ . Thus, it "moves" the value x' mod  $2^r$  to its "correct" location and we have to add  $f_a$  and subtract  $c_a$  afterward. However, implementations usually construct, instead of the arrays unshuffle and range, an array decode that stores, for each number x' mod  $2^r$ , an already corrected value for x and the corresponding letter a; hence, the decoding is much simpler:

$$(x, a) = \mathsf{decode}[x' \bmod 2^r].$$

The described scheme is called a *tabled ANS (tANS)*. It is the most popular variant of ANS widely used in practice. The choice of a shuffling method is crucial for its performance. Some methods are considered in the next section. There are several additional technical improvements that can be applied to this basic scheme. Perhaps, the most notable of them is that one can feed to the decoder more than r+1 bits at once, decoding many letters in one step (the information about the decoded letters and the new value for x must be stored in the array decode). We do not discuss these details further.

# 3.3 Shuffling Methods

The general rule for shuffling is to distribute letters in the array range as uniformly as possible so that, for any letter a, the distance between consecutive occurrences of a is approximately  $\frac{n}{f_a}$ . Implementations usually use heuristics for this [7] or Duda's method [12] (which is introduced below). Let us discuss some considerations on this regard that will be developed in a more rigorous way in Section 4 in the sequel.

The following informal argument justifies the scheme with shuffling. The value  $\log x' - \log x$  is, in a sense, an increase in bits from the number  $w_i$  to  $w_{i+1}$ ; the bit length of  $w_n$ , the final encoding, is approximately the sum of the increases (this reasoning is formalized in Section 4). Denote  $\delta = x \mod f_a$ . We have  $x = f_a + \delta$ . If the letters in the array range are distributed uniformly, the distance between two consecutive letters a is approximately  $\frac{n}{f_a}$ . Therefore, the encoder transforms x approximately to  $x' \approx 2^r + \frac{n}{f_a} \delta = \frac{n}{f_a} (f_a + \delta) = \frac{n}{f_a} x$  (recall that  $n = 2^r$ ). Hence, we obtain  $\log x' - \log x \approx \log \frac{n}{f_a}$ , which is precisely the optimal number of bits for the letter a according to the entropy formula.

The argument suggests that the shuffling method should spread the letters in the array range in such a way that the distance between two consecutive occurrences of a is approximately  $\frac{n}{f_a}$  and the encoder transforms the number  $x = f_a + \delta$  as close as possible to the number  $2^r + \frac{n}{f_a}\delta$ . Under this assumption of "uniformity", if the first occurrence of letter a is at position p in range, then x is transformed into  $x' \approx 2^r + \frac{n}{f_a}\delta + p$ . The term p adds to the redundancy associated with  $f_a$  letters a: the larger the value of p, the more bits are spent per letter a. Therefore, the first occurrences of more frequent letters should be closer to the beginning of the array range so that they produce less redundancies overall. Duda's method, which he called a precise initialization, tries to take into account all these considerations.

Duda's algorithm maintains a priority queue with the following operations: put(q, a) adds a pair of numbers (q, a) to the queue; (q, a) = getmin() removes from the queue a pair (q, a) with the smallest value q (breaking ties arbitrarily). We first give in Algorithm 1 simplified Duda's algorithm, which is easier to analyze.

**Algorithm 1** A simple initialization algorithm.

```
\begin{split} & \textbf{for } a \in [0..\sigma) \textbf{ do} \\ & \text{put}(0,a), \ d_a = c_a; \\ & \textbf{for } i = 0,1,\dots,2^r-1 \textbf{ do} \\ & (q,a) = \text{getmin}(), \ \text{range}[i] = a, \ \text{shuffle}[d_a] = i; \\ & \text{put}(q + \frac{n}{l_a},a), \ d_a = d_a + 1; \end{split}
```

The array range corresponding to shuffle is not needed for the encoder and its construction is added here for the convenience of the reader.

At every moment during the work of the algorithm, there is only one instance of each letter in the priority queue and, if a letter a was assigned to the array range exactly k times, then it is represented by the pair  $(\frac{n}{f_a}k,a)$  in the queue. Therefore, after  $f_a$  assignments of a into range, the letter is represented by (n,a) and all other letters b that had less than  $f_b$  assignments in range are represented as (q,b) with q < n. Hence, in the end, each letter a occurs in range exactly  $f_a$  times.

Duda's original "precise initialization" is the same as Algorithm 1 except that the operation "put(0, a)" from the first loop is changed to "put( $\frac{1}{2} \cdot \frac{n}{f_a}$ , a)". Its correctness is proved analogously.

# 4 Analysis of the Tabled ANS

We are to estimate the redundancy of the output produced by the ANS encoder, i.e., the difference between  $\lceil \log w_n \rceil$  and the information theoretic lower bound  $\sum_a f_a \log \frac{n}{f_a}$ . We postpone the analysis of the ANS without shuffling to the next section, where its more general variant is considered. In this subsection we consider the ANS with shuffling, i.e., the tabled ANS (tANS). To the best of our knowledge, the following analysis, albeit quite simple, evaded the attention of researchers and was not present in prior works. Our proof methods, however, stem from observations and arguments from Dubé and Yokoo [25] and Duda [11, 12].

# 4.1 Upper Bounds

Let us upperbound  $\log w_{i+1} - \log w_i = \log \frac{w_{i+1}}{w_i}$ , a bit increase after one step of the encoding procedure. The total number of bits will be then estimated as follows (the term r appears because  $w_0 = 2^r$  and  $\log w_0 = r$ ):

$$\log w_n = \log \frac{w_n}{w_{n-1}} + \log \frac{w_{n-1}}{w_{n-2}} + \dots + \log \frac{w_1}{w_0} + r.$$
 (6)

Suppose that  $w_{i+1}$  was obtained from  $w_i$  by "inserting" a letter a as described above. Denote  $\ell = \lfloor \log w_{i+1} \rfloor - r$  so that  $w_i = x \cdot 2^\ell + \Delta$  and  $w_{i+1} = x' \cdot 2^\ell + \Delta$ , where  $0 \le \Delta < 2^\ell$ . Then,  $\log \frac{w_{i+1}}{w_i} = \log \left(\frac{x'2^\ell + \Delta}{x^{2^\ell} + \Delta}\right) = \log \left(\frac{x'}{x}\left(\frac{1 + \Delta/(x'2^\ell)}{1 + \Delta/(x2^\ell)}\right)\right) = \log x' - \log x + \log \left(\frac{1 + \Delta/(x'2^\ell)}{1 + \Delta/(x2^\ell)}\right)$ . Since x' > x, the additive term  $\log \left(\frac{1 + \Delta/(x'2^\ell)}{1 + \Delta/(x2^\ell)}\right)$  is negative and, thus, we have obtained the following inequality:

$$\log w_{i+1} - \log w_i \le \log x' - \log x. \tag{7}$$

It remains to estimate how close is  $\log x' - \log x$  to the optimum  $\log \frac{n}{f_a}$ . We first consider the case when the encoder uses the shuffling produced by simplified Duda's algorithm (Algorithm 1).

Fix a letter a and a number  $\delta \in [0..f_a)$ . Denote by k the index of the  $(\delta + 1)$ th occurrence of a in range. Note that  $\mathsf{shuffle}[c_a + \delta] = k$ , by definition. For each  $b \in [0..\sigma)$ , denote by  $k_b$  the number of letters b in the subrange  $\mathsf{range}[0..k-1]$ . Clearly, we have  $k = \sum_b k_b$ . The shuffling algorithm implies the following inequality:

$$(k_b - 1)\frac{n}{f_b} \le \delta \frac{n}{f_a}. (8)$$

We express  $k_b$  from (8) as  $k_b \leq \delta \frac{f_b}{f_a} + 1$ . Summing over all  $b \in [0..\sigma)$ , we deduce from this  $k \leq \delta \frac{n}{f_a} + \sigma$ . It follows from (7) that, in order to analyze the number of bits per letter produced by the encoder, we have to estimate  $\log x' - \log x$ , where, by (4),  $x' = 2^r + \text{shuffle}[c_a + (x \mod f_a)]$ . Assuming  $\delta = x \mod f_a$ , we obtain  $x = f_a + \delta$  and  $x' = 2^r + k = n + k$ . Therefore,  $\log x' - \log x = \log(n + k) - \log x \leq \log(n + \delta \frac{n}{f_a} + \sigma) - \log x = \log(\frac{n}{f_a} + \sigma) - \log x = \log(\frac{n}{f_a} + \log(1 + \frac{\sigma}{nx/f_a})) \leq \log(\frac{n}{f_a} + \frac{\sigma}{n}) \log e$ . Thus, we estimate the number of bits per letter a as  $\log(\frac{n}{f_a} + \frac{\sigma}{n}) \log e$ , i.e., the redundancy is  $\frac{\sigma}{n} \log e$  bits per letters.

Now let us analyze Duda's original algorithm. The algorithm is the same as Algorithm 1 except that the operation "put(0, a)" from the first loop is changed to "put( $\frac{1}{2} \cdot \frac{n}{f_a}$ , a)". The analysis is slightly more complicated. First, an equation analogous to (8) for this case looks as follows:

$$\left(k_b - \frac{1}{2}\right) \frac{n}{f_b} \le \left(\delta + \frac{1}{2}\right) \frac{n}{f_a}.$$

We then similarly deduce  $k_b \leq (\delta + \frac{1}{2}) \frac{f_b}{f_a} + \frac{1}{2}$  and, summing over all  $b, k \leq (\delta + \frac{1}{2}) \frac{n}{f_a} + \frac{\sigma}{2}$ . Again, assuming  $x = f_a + \delta$ , it follows from this that  $\log x' - \log x \leq \log(\frac{n}{f_a}x + \frac{n}{2f_a} + \frac{\sigma}{2}) - \log x = \log\frac{n}{f_a} + \log(1 + \frac{1}{2x} + \frac{\sigma}{2nx/f_a}) \leq \log\frac{n}{f_a} + (\frac{1}{2f_a} + \frac{\sigma}{2n})\log e$ . The letter a occurs in the sequence exactly  $f_a$  times. Hence, the redundancies  $\frac{1}{2f_a}\log e$  for letters a sum up to  $\frac{1}{2}\log e$  over all these occurrence and, therefore, in the end we obtain  $\frac{\sigma}{2n}\log e$  bits per letter contributed by the terms  $\frac{1}{2f_a}\log e$ , for all letters a, in the final encoding. Adding to this the  $\frac{\sigma}{2n}\log e$  bits per letter, we obtain  $\frac{\sigma}{n}\log e$  redundant bits per letter in the final encoding.

It remains to add the additive term r from (6) to the redundancy, which contributes  $\frac{r}{n}$  bits per letter, and the following theorem is proved.

**Theorem 1.** Given a sequence  $a_1, a_2, \ldots, a_n$  of letters from an alphabet  $[0..\sigma)$  such that each letter a occurs in it  $f_a$  times and  $n = 2^r$  for an integer r, the ANS encoder using [simplified] Duda's precise initialization transforms this sequence into a bit string of length

$$\sum_{a \in [0..\sigma)} f_a \cdot \log \frac{n}{f_a} + O(\sigma + r),$$

where the  $O(\sigma + r)$  term can be bounded by  $\sigma \log e + r$ . Thus, we have  $O(\frac{\sigma + r}{n})$  redundant bits per letter.

In practice, the length m of the encoded sequence  $a_1, a_2, \ldots, a_m$  is not necessarily a power of two. A typical solution for this case is to approximate the real empirical probabilities  $f_a/m$  of letters with approximate ones  $\hat{f}_a/2^r$ , where  $\sum_a \hat{f}_a = 2^r$ . The ANS encoder then processes the sequence as usually but using the "frequencies"  $\hat{f}_a$  instead of  $f_a$ . The same analysis can be applied for this case: Equations (6) and (7) trivially hold and the value  $\log x' - \log x$  is bounded in the same manner by  $\log \frac{2^r}{\hat{f}_a} + \frac{\sigma}{2^r} \log e$ , for simplified Duda's initialization, and by  $\log \frac{2^r}{\hat{f}_a} + (\frac{1}{2\hat{f}_a} + \frac{\sigma}{2 \cdot 2^r}) \log e$ , for Duda's initialization. Summing the redundancies over all m letters, we obtain the following theorem (for simplicity, the theorem is stated only for simplified Duda's initialization).

**Theorem 2.** Let  $a_1, a_2, \ldots, a_m$  be a sequence of letters from an alphabet  $[0..\sigma)$  such that each letter a occurs in it  $f_a$  times. Let the probabilities  $f_a/m$  be approximated by numbers  $\hat{f}_a/n$  such that  $\hat{f}_a$  are integers,  $n=2^r$ , for an integer r, and  $\sum_{a\in[0..\sigma)}\hat{f}_a=n$ . The ANS encoder that uses the approximate probabilities and simplified Duda's initialization transforms this sequence into a bit string of length

$$\sum_{a \in [0..\sigma)} f_a \cdot \log \frac{n}{\hat{f}_a} + O(\frac{\sigma m}{n} + r),$$

where the  $O(\frac{\sigma m}{n} + r)$  term can be bounded by  $\frac{\sigma m}{n} \log e + r$ .

### 4.2 Lower Bound Example

Apparently, the r-bit redundancy incurred by the initial value  $w_0$  is unavoidable in the described scheme. It is less clear whether an  $O(\sigma)$  additive term is necessary in Theorem 1. An informal argument supporting that this is the case is as follows. Consider a sequence in which all letters are (approximately) equiprobable, i.e., their frequencies  $f_a$  are  $\sim \frac{n}{\sigma}$ . The lower bound for the encoding of the sequence is  $n\log\sigma$  bits. The array range constructed by Duda's initialization algorithm for the sequence looks (approximately) as  $n/\sigma$  blocks, each of which is of size  $\sigma$  and consists of consecutive letters  $0,1,\ldots,\sigma-1$ . Hence, when the encoder receives a letter  $a_{i+1}=a\in[0..\sigma)$  during its work, it transforms a number  $x=f_a+\delta$ , where  $\delta=x$  mod  $f_a$ , occupying leading bits of  $w_i$ , into the number  $x'=2^r+\frac{n}{f_a}\delta+a=\frac{n}{f_a}x+a=\sigma x+a$  (note that  $f_a=n/\sigma$ , by our assumption). As in the previous section, one can deduce from this that  $\log x'-\log x=\log\sigma+\log(1+\frac{a}{\sigma x})$ . Since  $x<2f_a=2\frac{n}{\sigma}$ , the redundant additive term  $\log(1+\frac{a}{\sigma x})$  can be estimated as  $\log(1+\frac{a}{\sigma x})\geq\log(1+\frac{a}{2n})\geq\frac{a}{2n}\cdot\frac{1}{1+a/(2n)}$ . Thus, the redundancy is approximately  $\frac{a}{2n}$  bits per letter a, which sums to the total redundancy of  $\sum_a \frac{a}{2n}f_a=\sum_a \frac{a}{2\sigma}=\frac{\sigma+1}{4}$  bits over all letters in the sequence.

This informal argument is only an intuition since the negative terms  $\log\left(\frac{1+\Delta/(x'2^{\ell})}{1+\Delta/(x2^{\ell})}\right)$  that appear in the analysis of Section 4.1 could, in principle, diminish the described effect. Nevertheless, as we are to show, an  $\Omega(\sigma)$  redundancy indeed appears in some instances.

Fix an even integer r > 0. Observe that  $2^r \equiv 1 \pmod{3}$  since r is even. Denote  $n = 2^r$ . The sequence under construction will contain  $\sigma = (n-1)/3+1$  letters  $0,1,\ldots,\sigma-1$ . Each letter  $a \in [0..\sigma-1)$  has exactly three occurrences in the sequence (i.e.,  $f_a = 3$ ) and the letter  $\sigma - 1$  occurs only once (i.e.,  $f_{\sigma-1} = 1$ ); note that  $\sum_{a \in [0..\sigma)} f_a = 3(\sigma - 1) + 1 = n$ . The entropy formula gives the following lower bound on the encoding size for the sequence:

$$(n-1)\log\frac{n}{3} + \log n = (n-1)(r-\log 3) + r = (n-1)(r-1.58496...) + r.$$
(9)

It is straightforward that with such frequencies of letters both Duda's initialization algorithm and its simplified variant (Algorithm 1) construct the same array range: the subrange range $[0..\sigma-1]$  contains consecutively the letters  $0,1,\ldots,\sigma-1$  (in this order) and the subranges range $[\sigma..2\sigma-2]$  and range $[2\sigma-1..n-1]$  are equal and both contain consecutively the letters  $0,1,\ldots,\sigma-2$  (in this order). Now let us arrange the letters in the sequence  $a_1,a_2,\ldots,a_n$ .

The last letter  $a_n$  is  $\sigma-1$ . The rest,  $a_1, a_2, \ldots, a_{n-1}$ , consists of letters  $a \in [0..\sigma-1)$  whose frequencies are  $f_a=3$  (11<sub>2</sub> in binary). When the encoder processes a letter  $a_{i+1}=a \in [0..\sigma-1)$  and modifies the number  $w_i$  representing the prefix  $a_1, a_2, \ldots, a_i$ , it replaces either two or three leading bits of  $w_i$  with new r+1 bits, thus producing the number  $w_{i+1}$ . The choice of whether to replace two or three bits depends on whether the two leading bits of  $w_i$  are 11 or 10, respectively (i.e., whether the two bits store a number less than  $f_a=3$  or not). We are to arrange the letters  $[0..\sigma-1)$  in the sequence  $a_1, a_2, \ldots, a_{n-1}$  in such a way that the encoder chooses the two options alternatingly: it replaces three leading bits of  $w_i$  if i is even, and two bits if i is odd. The total number of bits produced in this way is at least  $\frac{n-1}{2}(r-1) + \frac{n-1}{2}(r-2) + r = (n-1)(r-1.5) + r$  (the additive term r appears when the last letter  $a_n = \sigma - 1$  is encoded). Comparing this to (9), we can see that the encoding generated by ANS is larger than the optimum (9) by at least  $(\log 3 - 1.5)(n-1) > 0.08496(n-1)$ . By simple calculations, we deduce from the equality  $\sigma = (n-1)/3 + 1$  that the redundancy 0.08496(n-1) is larger than  $\frac{\sigma-1}{4}$ . Let us describe an arrangement of letters that produces such effect of "alternation".

The encoder consecutively transforms the initial value  $w_0 = 2^r$  into  $w_1, w_2, \ldots$  by performing the push operations:  $w_{i+1} = \mathsf{push}(w_i, a_{i+1})$ . Let us call the number  $\lfloor w_i/2^{\lfloor \log w_i \rfloor - r} \rfloor \mod 2^r$  a state; it is the value stored in the highest r+1 bits of the number  $w_i$  currently processed by the encoder minus the highest bit 1. The range of possible states is  $[0..2^r) = [0..n)$ . We split this range into three disjoint segments numbered i, ii, iii:  $[0..\sigma)$ ,  $[\sigma..2\sigma - 1)$ ,  $[2\sigma - 1..n)$ , whose lengths are  $\sigma$ ,  $\sigma - 1$ ,  $\sigma - 1$ , respectively. The subarray range  $[0..\sigma-1]$  corresponding to the segment i contains all letters  $[0..\sigma)$ ; each of the subarrays range  $[\sigma..2\sigma - 2]$  and range  $[2\sigma - 1..n - 1]$  corresponding to the segments ii and iii, respectively, contains all letters  $[0..\sigma-1)$ .

Receiving the letter  $\sigma - 1$  (which occurs in the sequence only once), the encoder transforms any current state to the state  $\sigma - 1$  (it is the index of the letter  $\sigma - 1$  in the array range). In order to describe how

states are transformed by other letters, let us split the state range  $[0..2^r)$  into three segments called A, B, C:  $[0..2^{r-2})$ ,  $[2^{r-2}..2^{r-1})$ ,  $[2^{r-1}..n)$ , whose lengths are  $2^{r-2}$ ,  $2^{r-2}$ ,  $2^{r-1}$ , respectively. The segments are related as follows: receiving a letter from  $[0..\sigma-1)$ , the encoder translates the current state from the segment A (respectively, B, C) to a state from the segment ii (respectively, iii, i); see Figure 5 where this transition of states is illustrated by dashed lines connecting the corresponding segments.

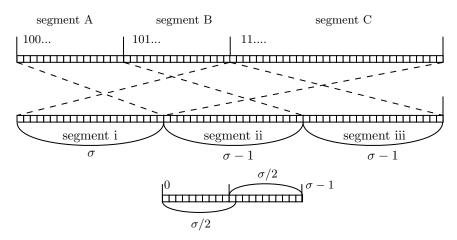


Figure 5: Two splits of the state range  $[0..2^r)$  into segments. Here r=6. Common first bits of numbers  $w_i$  corresponding to the states belonging to the segments A, B, C are written above the respective segments. Observe that the bits start with 11 only when a state belongs to the segment C. The range  $[0..\sigma-1)$  is depicted under the segment ii: the encoder receives letters from the right part  $[\sigma/2 - 1..\sigma - 1)$  when the current state is from the segment A, and from the left part  $[0..\sigma/2)$  when the state is from the segment C.

Receiving  $a_{i+1} = a \in [0..\sigma-1)$ , the encoder replaces two leading bits of  $w_i$  with r+1 new bits if the current state is from the segment C; otherwise, it replaces three leading bits of  $w_i$ . This behaviour is determined by first bits of  $w_i$ , which are equal to 11... only for states from the segment C; see Figure 5. The execution of the encoder starts with the state 0, which belongs to the segment A. We are to arrange letters in the sequence  $a_1, a_2, \ldots, a_{n-1}$  so that, for odd i, the state corresponding to  $w_i$  belongs to the segment C, and for even i, to the segment A. The states in our arrangement will never belong to the segment B.

Since  $\sigma + \sigma/2 - 1 = n/2$  and n/2 is the leftmost state from the segment C, any state from the segment A transits to a state from the segment C when the encoder receives a letter a from the range  $[\sigma/2 - 1..\sigma - 1)$ ; the new state is  $\sigma + a$ , the ath element of the segment ii (see Figure 5 for an illustration). Similarly, since  $\sigma/2 < n/4$  and n/4 - 1 is the rightmost state from the segment A, any state from the segment C transits to a state from the segment A when the encoder receives a letter  $a \in [0..\sigma/2)$ . Accordingly, for  $i \in [1..n)$ , we put in the sequence as the letter  $a_i$  a letter from  $[\sigma/2 - 1..\sigma - 1)$  if i is odd, and a letter from  $[0..\sigma/2)$  if i is even (note that both ranges share a common letter  $\sigma/2 - 1$ ; it is not a mistake). The states thus "bounce" between the segments A and C during the processing of  $a_1, a_2, \ldots, a_{n-1}$  by the encoder. Since the sizes of both ranges  $[0..\sigma/2)$  and  $[\sigma/2 - 1..\sigma - 1)$  are  $\sigma/2$  and they share a common letter  $\sigma/2 - 1$ , the letters  $[0..\sigma - 1)$  can be distributed in the sequence  $a_1, a_2, \ldots, a_{n-1}$  in such way that each letter occurs exactly three times and each letter  $a_i$  is from the range  $[\sigma/2 - 1..\sigma - 1)$ , for odd i, and from the range  $[0..\sigma/2)$ , for even i.

and each letter  $a_i$  is from the range  $[\sigma/2 - 1..\sigma - 1)$ , for odd i, and from the range  $[0..\sigma/2)$ , for even i. We thus have obtained a redundancy of  $(\log 3 - 1.5)(n-1) > \frac{\sigma-1}{4}$  bits. Adding to this r-2 bits produced by the r-2 lowest bits of the initial value  $w_0$ , we have proved the following theorem.

**Theorem 3.** For arbitrarily large  $n = 2^r$ , there exists a sequence  $a_1, a_2, \ldots, a_n$  of letters from an alphabet  $[0..\sigma)$  with  $\sigma > n/3$  such that the ANS encoder using [simplified] Duda's precise initialization transforms this sequence into a bit string of length at least

$$\sum_{a \in [0..\sigma)} f_a \cdot \log \frac{n}{f_a} + \frac{\sigma - 1}{4} + r - 2,$$

where  $f_a$  is the number of occurrences for letter a. Thus, the redundancy is  $\frac{\sigma-1}{4}+r-2$  bits.

The example that attains the lower bound of Theorem 3 is very simple. Hence, it is reasonable to assume that any adequate shuffling method would have the same redundancy  $\Omega(\sigma+r)$  as in Theorem 3 on a similarly constructed sequence  $a_1, a_2, \ldots, a_n$ . We believe, therefore, that Duda's conjecture that the redundancy can be  $O(\frac{\sigma}{r^2})$  bits per letter when an appropriate shuffling method is used is disproved.

# 5 Range ANS with Fixed Accuracy

In Section 3.1, a simple range ANS (rANS) was described, which is just the ANS without shuffling. Now we are to introduce another variant of rANS, called rANS with fixed accuracy to distinguish it from the rANS as defined by Duda [12, 14] (which is also sketched below). Our exposition is less detailed than in the previous sections since we believe that all ideas and intuition necessary for understanding were developed above.

Let  $a_1, a_2, \ldots, a_n$  be a sequence of letters over an alphabet  $[0..\sigma)$ , where  $n=2^r$  and the frequencies of letters are denoted by  $f_a$  (i.e., the probability for letter a is  $\frac{f_a}{n}$ ). Fix an integer  $k \geq 0$ , which will serve as a user-defined accuracy parameter regulating the size of redundancy (typically,  $0 \leq k \leq 4$ ). As in the simple rANS, the encoder starts its work with a number  $w_0 = 2^{r+k}$ ; the value  $w_0$  might be arbitrary, the only necessary condition is  $w_0 \geq 2^{r+k}$ . The encoder consecutively performs operations  $w_{i+1} = \text{push}(w_i, a_{i+1})$ , for  $i = 0, 1, \ldots, n-1$ , but the operation push(w, a) works differently this time: it substitutes either  $\lfloor \log f_a \rfloor + k + 1$  or  $\lfloor \log f_a \rfloor + k + 2$  highest bits of w that store a number x by x + k + 1 new bits that store the number  $x' = \lfloor x/f_a \rfloor 2^r + c_a + (x \mod f_a)$ , where  $x' = \sum_{b \in [0..a)} f_b$ ; the condition determining the number of bits occupied by x is essentially the same as in the simple rANS from Section 3.1.

More formally, the algorithm for  $\operatorname{push}(w,a)$  is as follows. Denote by  $x_1$  and  $x_2$  the values stored in, respectively, the highest  $\lfloor \log f_a \rfloor + k + 1$  and  $\lfloor \log f_a \rfloor + k + 2$  bits of w. We assume that  $x = x_1$  if  $x_1 \geq f_a 2^k$ , and  $x = x_2$  otherwise. Since  $2^{\lfloor \log f_a \rfloor + k} \leq x_1 < f_a 2^{k+1}$  and  $x_2 \leq 2 \cdot x_1 + 1$ , we have  $f_a 2^k \leq x < f_a 2^{k+1}$ . Denote  $\ell = \lfloor \log w \rfloor - \lfloor \log x \rfloor$ . Note that  $w = x \cdot 2^{\ell} + (w \mod 2^{\ell})$ . The value  $w' = \operatorname{push}(w,a)$  is computed as  $w' = x' \cdot 2^{\ell} + (w \mod 2^{\ell})$  by replacing the part x of w by x + k + 1 bits representing a number x' defined as:

$$x' = \lfloor x/f_a \rfloor 2^r + c_a + (x \bmod f_a). \tag{10}$$

Since  $f_a 2^k \le x < f_a 2^{k+1}$ , we have  $2^k \le \lfloor x/f_a \rfloor < 2^{k+1}$  and, therefore, the number x' indeed fits into r+k+1 bits and its highest (r+k+1)th bit is 1. The reverse operation  $(w,a) = \mathsf{pop}(w')$  producing the old value w and the letter a from w' is straightforward. Given a number x', which occupies r+k+1 highest bits of w' (i.e.,  $x' = \lfloor w'/2^\ell \rfloor$ , where  $\ell = \lfloor \log w' \rfloor - r - k$ ), we first determine the letter a by examining to which range  $\lfloor c_a ... c_{a+1} \rfloor$  the number  $x' \mod 2^r$  belongs and, then, we compute x as follows:

$$x = \lfloor x'/2^r \rfloor f_a + (x' \bmod 2^r) - c_a. \tag{11}$$

Once x is known, we put  $w = x \cdot 2^{\ell} + (w' \mod 2^{\ell})$ . It remains to observe that x' > x since  $x' = \lfloor x/f_a \rfloor 2^r + (x \mod f_a) > \lfloor x/f_a \rfloor f_a + (x \mod f_a) = x$ . Therefore, we have  $w_0 < w_1 < \cdots < w_n$ .

Note that, when k = 0, the described scheme degenerates simply to the ANS without shuffling.

To estimate the size of the final number  $w_n$  in bits, we derive by analogy to (6) and (7) using the condition x' > x the following two equations (here we have  $w_i = x \cdot 2^{\ell} + \Delta$  and  $w_{i+1} = x' \cdot 2^{\ell} + \Delta$ , where  $\Delta \in [0...2^{\ell})$ ):

$$\log w_n = \log \frac{w_n}{w_{n-1}} + \log \frac{w_{n-1}}{w_{n-2}} + \dots + \log \frac{w_1}{w_0} + r + k;$$

$$\log w_{i+1} - \log w_i \le \log x' - \log x.$$

From (11), we deduce  $\log x \ge \log(\lfloor x'/2^r \rfloor f_a) \ge \log f_a + \log(x'/2^r - 1) = \log f_a + \log(x'/2^r) + \log(1 - 2^r/x') = \log \frac{f_a}{2^r} + \log x' + \log(1 - 2^r/x')$ . Since  $2^{r+k} \le x'$ , we derive further  $\log(1 - 2^r/x') \ge \log(1 - 1/2^k) \ge -\frac{\log e}{2^k - 1}$ . Therefore, we obtain

$$\log x' - \log x \le \log \frac{n}{f_a} + \frac{\log e}{2^k - 1}.$$

Summing the values  $\log x' - \log x$  over all  $n = 2^r$  letters of the sequence, we obtain the following theorem.

**Theorem 4.** Given a sequence  $a_1, a_2, \ldots, a_n$  of letters from an alphabet  $[0..\sigma)$  such that each letter a occurs in it  $f_a$  times and  $n = 2^r$  for an integer r, the rANS encoder with fixed accuracy  $k \ge 1$  transforms this sequence into a bit string of length

$$\sum_{a \in [0..\sigma)} f_a \cdot \log \frac{n}{f_a} + O\left(\frac{n}{2^k} + r\right),\,$$

where the  $O\left(\frac{n}{2^k}+r\right)$  redundancy term can be bounded by  $\frac{n \log e}{2^k-1}+r$ .

For completeness, let us briefly describe the standard rANS [12, 14]. The encoder similarly starts with the number  $w_0 = 2^r$  and consecutively computes  $w_1, w_2, \ldots, w_n$  for the sequence  $a_1, a_2, \ldots, a_n$ , where  $n = 2^r$ . The encoder maintains a number t (initially t = r + 1) and, receiving a new letter, it replaces the highest t bits of the current number  $w_i$  with a larger value and increases t accordingly. A "renormalization" is sometimes performed by reducing the value t in order to contain numbers within a range fitting into a machine word. More precisely, receiving a letter  $a = a_{i+1}$ , the encoder takes the value t stored in the t highest bits of t (i.e., t with t with t thus producing t with t thus producing t the encoder takes the value t is increased by t formula (10), and replaces t with t thus producing t thus producing t the encoder takes the value t is increased by t formula (10), and replaces t with t thus producing t thus producing t the encoder takes the value t is increased by t formula (10), and replaces t with t thus producing t thus producing t the encoder takes the value t is increased by t formula (10), and replaces t with t thus producing t t

The analysis of this rANS variant is not in the scope of the present paper; see [12, 22].

Implementation notes In our experiments the described rANS encoder with fixed accuracy k = 3 was 1.8 times slower than the tANS encoder implemented by Collet [7] and had approximately the same compression rate.<sup>2</sup> It is unsurprising since the inner encoding loop of tANS computing (2) essentially consists of just a couple of accesses to tables stored in the L1 cache. However, the rANS with fixed accuracy is faster than one could have expected from the standard rANS. The key feature that allows the speed boost is that the operation of division  $\lfloor x/f_a \rfloor$  in (10) guarantees that its result is in the range  $\lfloor 2^k ... 2^{k+1} \rfloor$ . Therefore, the division can be executed by simpler instructions in a branchless code: we used arithmetic and bit operations and the instruction cmov from x86 (it is also possible to use only arithmetic and bit instructions). The code, however, turns out to be quite cumbersome, which noticeably diminishes the benefits of the division-free branchless loop.

Denote R = r + k. The main loop of the encoder calls the function encode from Algorithm 2 consecutively for the letters  $a_1, a_2, \ldots, a_n$  in the encoded sequence. The function receives as its parameters an (R+1)-bit number w and a letter a. The function stores some lowest bits of w in an external storage and returns an (R+1)-bit value x' computed as in (10) (details follow). The parameter w is actually an (R+1)-bit number x' produced by the previous call to the function encode in the encoding loop; the first call receives  $w = 2^R$ .

Denote  $t = r - \lfloor \log f_a \rfloor$ . The number x occupies either R - t + 1 or R - t + 2 highest bits of w. The presented pseudocode uses Collet's trick [7] to determine x with a branchless code. To this end, the array table stores, for each letter  $a \in [0..\sigma)$ , besides the values  $c_a$  and  $f_a$  the number  $d = (t << (R+1)) - (f_a << (t+k))$ . The trick is that, in this case, the number  $w + d = (t << (R+1)) + w - (f_a << (t+k))$  contains in its highest bits  $R+1, R+2, \ldots$  either the number t or t-1 depending on whether  $w \ge (f_a << (t+k))$  or not. Therefore, x = w >> s, where s = (w+d) >> (R+1).

The code in lines 6–11 accumulates the quotient  $\lfloor x/f_a \rfloor$  in the variable q and the remainder  $x \mod f_a$  in the variable x. It is done by subtracting the numbers  $f_a << i$ , for i=3,2,1,0, from x, thus, reconstructing q bit by bit; note, however, that the bits in q are inverted and, hence, in the end we have to perform xor with 15 (1111<sub>2</sub> in binary).

<sup>&</sup>lt;sup>2</sup>In [16] Giesen showed that interleaving streams can significantly speed up the rANS. We did not use this trick, however, because Collet's tANS does not use it either, though we believe that tANS can utilize this technique too, albeit less efficiently.

#### **Algorithm 2** The encoding function of rANS with fixed accuracy k = 3.

```
\triangleright rANS with fixed accuracy k=3
 1: function encode(w, a)
                                                    \Rightarrow d = (t << (R+1)) - (f_a << (t+k)), \text{ where } t = r - \lfloor \log f_a \rfloor
        (c_a, f_a, d) = \mathsf{table}[a];
        s = (w+d) >> (R+1);
 3:
                                                                  \triangleright output s lowest bits of w into an external storage
        outputBits(w, s);
 4:
        x = w >> s;
 5:
        x = x - (f_a << 3);
 6:
 7:
        q = 0;
        for i = 2, 1, 0 \text{ do}
                                                                                               ▶ the loop must be unrolled
 8:
            x_0 = x - (f_a << i);
9:
            q = q or (x_0 and (1 << (R+i)));
10:
            if (x_0 \ge 0) \ x = x_0;
                                                                                         ▶ it is compiled into cmov on x86
11:
        return ((q \text{ xor } (15 << R)) >> k) + c_a + x;
12:
```

The standard rANS has a couple of advantages over tANS in some use cases [16]: it does not require a table of size  $2^r$  like tANS and, hence, is more convenient for the (pseudo) adaptive mode when the table of frequencies is sometimes rebuilt during the execution of the encoder; due to this less heavy use of memory, the rANS might be better for interleaving several streams of data and utilizes more efficiently the instruction-level parallelism for this task. For these reasons, the rANS was used in some high performance compressors [16]. The described rANS with fixed accuracy shares the same good features of the standard rANS plus the described above benefits of the controlled division. In addition, we believe that the rANS with fixed accuracy can potentially have more efficient hardware implementations using a parallelization for the computation of the division, which is possible since the resulting quotient is in a small range  $[2^k...2^{k+1})$ .

# 6 Conclusion and Open Problems

Theorems 1 and 3 describe the tight asymptotic behaviour of the redundancy for the tANS. We believe that Theorem 4, albeit not complemented with a lower bound, is asymptotically tight too. However, it is open to provided a series of examples supporting this claim. A number of other problems listed below still remain open too.

- (i) The main remaining open problem concerning ANS encoders, as we see it, is to construct a genuine FIFO encoder with the same performance characteristics as tANS or rANS. The known ANS variants work as stacks (LIFO) while it is more natural and, in some scenarios, preferable to have an encoder that acts as a queue, like the arithmetic coding that fulfils this requirement but is noticeably worse than ANS in performance terms. For the same reason, the current ANS variants are not suitable enough for the adaptive encoding when frequencies of letters change as the algorithm reads the sequence from left to right (however, the ANS can be applied for a pseudoadaptive mode when the input sequence is split into chunks and the frequencies change after a whole chunk is processed; the rANS is primarily used for this purpose since it requires less tables to rebuild [16]).
- (ii) Duda's precise initialization and its simplified variant from Algorithm 1 are not particularly suitable for practice due to the overhead incurred by operations on the priority queue and by floating point operations. Therefore, usually they are replaced with heuristics. An implementation of a fast and simple initialization method with good guarantees sufficient for Theorem 1 is an open problem.
- (iii) The constants in our lower and upper bounds in Theorems 1 and 3 do not coincide and it remains open to find a tight constant for the tANS redundancy term. We believe also that the r-bit redundancy produced by the initial state  $w_0$  of the encoder can be somehow reduced too by slightly modifying the scheme. Clearly, one can omit the trailing zeros in the bit representation of the resulting number  $w_n$  but it is not enough to get rid of the r-bit redundancy entirely.
  - (iv) Another problem that can be addressed is the encoding of the table of frequencies (see [13]). We are

not aware of a suitable lower bound like the entropy formula that would include both the encoded sequence and the table of frequencies. We believe that such a formula could be derived by analyzing logarithms of multinomial coefficients and it is likely that the formula will confirm that the current simple methods for encoding the frequency table are close to optimal.

- (v) We believe that the implementation of the rANS with fixed accuracy presented in Algorithm 2 is not satisfactory and can be improved. The issue is that it uses too many instructions and the cmov instructions, in particular, which are relatively slow. A solution could be in precomputing some numbers (akin to numbers d in Collet's trick) that would help to speed up the computations.
- (vi) In practice, for performance reasons, the empirical probabilities of letters in the input sequence are usually approximated with probabilities of the form  $f/2^r$ . Due to Theorem 2, the redundancy incurred by these approximations can be (at least partially) estimated using the Kullback-Leibler distance. The problem of the choice for the approximations in practice is not so trivial and can be studies further; see [26].

## References

- [1] Apple LZFSE compressor, https://github.com/lzfse/lzfse, accessed: 07.01.2022.
- [2] J. Alakuijala, R. van Asseldonk, S. Boukortt, M. Bruse, I.-M. Comşa, M. Firsching, T. Fischbacher, E. Kliuchnikov, S. Gomez, R. Obryk, et al., JPEG XL next-generation image compression architecture and coding tools, in: Applications of Digital Image Processing XLII, vol. 11137, International Society for Optics and Photonics, 2019.
- [3] T. Alonso, G. Sutter, J. E. L. De Vergara, LOCO-ANS: An optimization of JPEG-LS using an efficient and low-complexity coder based on ANS, IEEE Access 9 (2021) 106606–106626.
- [4] I. Blanes, M. Hernandez-Cabronero, J. Serra-Sagrista, M. W. Marcellin, Redundancy and optimization of tANS entropy encoders, IEEE Transactions on Multimedia (2020) 4341–4350.
- [5] C. Bloom, cbloom rants blog, http://cbloomrants.blogspot.com, accessed: 08.01.2022.
- [6] Y. Collet, Facebook Zstandard compressor, https://github.com/facebook/zstd, accessed: 07.01.2022.
- [7] Y. Collet, FSE implementation, https://github.com/Cyan4973/FiniteStateEntropy, accessed: 07.01.2022.
- [8] Y. Collet, RealTime Data Compression blog, http://fastcompression.blogspot.com, accessed: 08.01.2022.
- [9] T. M. Cover, J. A. Thomas, Information theory and statistics, Elements of Information Theory 1 (1) (1991) 279–335.
- [10] D. Dubé, H. Yokoo, Fast construction of almost optimal symbol distributions for asymmetric numeral systems, in: 2019 IEEE International Symposium on Information Theory (ISIT), IEEE, 2019.
- [11] J. Duda, Asymmetric numeral systems, arXiv preprint arXiv:0902.0271 (2009) 1-47.
- [12] J. Duda, Asymmetric numeral systems: entropy coding combining speed of huffman coding with compression rate of arithmetic coding, arXiv preprint arXiv:1311.2540 (2013) 1–24.
- [13] J. Duda, Encoding of probability distributions for asymmetric numeral systems, arXiv preprint arXiv:2106.06438 (2021) 1–5.
- [14] J. Duda, K. Tahboub, N. J. Gadgil, E. J. Delp, The use of asymmetric numeral systems as an accurate replacement for Huffman coding, in: 2015 Picture Coding Symposium (PCS 2015), IEEE, 2015.

- [15] F. Giesen, The ryg blog, https://fgiesen.wordpress.com, accessed: 07.01.2022.
- [16] F. Giesen, Interleaved entropy coders, arXiv preprint arXiv:1402.3392 (2014) 1–16.
- [17] D. A. Huffman, A method for the construction of minimum-redundancy codes, Proc. Institute of Radio Engineers (IRE) 40 (9) (1952) 1098–1101.
- [18] G. N. N. Martin, Range encoding: an algorithm for removing redundancy from a digitised message, in: Proc. Institution of Electronic and Radio Engineers International Conference on Video and Data Recording, 1979.
- [19] A. Moffat, R. M. Neal, I. H. Witten, Arithmetic coding revisited, ACM Transactions on Information Systems 16 (3) (1998) 256–294.
- [20] A. Moffat, M. Petri, Large-alphabet semi-static entropy coding via asymmetric numeral systems, ACM Transactions on Information Systems 38 (4) (2020) 1–33.
- [21] J. J. Rissanen, Generalized Kraft inequality and arithmetic coding, IBM Journal of research and development 20 (3) (1976) 198–203.
- [22] J. Townsend, A tutorial on the range variant of asymmetric numeral systems, arXiv preprint arXiv:2001.09186 (2020) 1–11.
- [23] N. Wang, C. Wang, S.-J. Lin, A simplified variant of tabled asymmetric numeral systems with a smaller look-up table, Distributed and Parallel Databases 39 (3) (2021) 711–732.
- [24] H. Yokoo, On the stationary distribution of asymmetric binary systems, in: 2016 IEEE International Symposium on Information Theory (ISIT), IEEE, 2016.
- [25] H. Yokoo, D. Dubé, Asymptotic optimality of asymmetric numeral systems, in: Proc. 42nd International Symposium on Information Theory and Its Applications (SITA 2019, 2019.
- [26] H. Yokoo, T. Shimizu, Probability approximation in asymmetric numeral systems, in: Proc. 2018 International Symposium on Information Theory and Its Applications (ISITA 2018), IEEE, 2018.