

HANDBOOK ON ALGORITHMS AND THEORY OF
COMPUTATION

Text data compression algorithms

Contents

1	Text data compression algorithms	5
1.1	Text compression	5
1.2	Static Huffman coding	7
1.2.1	Encoding	8
1.2.2	Decoding	12
1.3	Dynamic Huffman coding	14
1.3.1	Encoding	15
1.3.2	Decoding	18
1.3.3	Updating	19
1.4	Arithmetic coding	22
1.4.1	Encoding	22
1.4.2	Decoding	23
1.4.3	Implementation	24
1.5	LZW Coding	32
1.5.1	Encoding	32
1.5.2	Decoding	33
1.5.3	Implementation	34
1.6	Mixing several methods	37
1.6.1	Run Length Encoding	37
1.6.2	Move To Front	38

1.6.3	Integrated example	38
1.7	Experimental results	40
1.8	Research Issues and Summary	41
1.9	Defining Terms	42
1.10	References	42
1.11	Further Information	43

List of Figures

1.1	Counts the character frequencies.	8
1.2	Builds the Huffman coding tree.	9
1.3	Builds character codewords from the coding tree.	10
1.4	Stores the coding tree in the compressed file.	10
1.5	Encodes the characters in the compressed file.	11
1.6	Complete function for Huffman encoding.	11
1.7	Rebuilds the tree from the header of compressed file.	13
1.8	Recovers the original text.	13
1.9	Complete function for Huffman decoding.	14
1.10	Initializes the dynamic Huffman tree.	16
1.11	Complete function for dynamic Huffman encoding.	16
1.12	Encodes one symbol.	17
1.13	Adds a new symbol in the tree.	17
1.14	Complete function for dynamic Huffman decoding.	18
1.15	Decodes one symbol from the input file.	19
1.16	Updates the current Huffman tree.	19
1.17	Basic arithmetic encoding.	22
1.18	Basic arithmetic decoding.	23
1.19	Complete arithmetic encoding function.	25
1.20	Encodes one symbol.	25

1.21	Sends one bit followed by the waiting bits, if any.	25
1.22	Complete arithmetic decoding function.	28
1.23	Decodes one symbol.	29
1.24	LZW encoding algorithm.	35
1.25	LZW decoding algorithm.	36
1.26	Example of text $y = \mathbf{baccara}$. Top line is $BW(y)$ and bottom line the sorted list of letters of it. Top-down arrows correspond to succession of occurrences in y . Each bottom-up arrow links the same occurrence of a letter in y . Arrows starting from equal letters do not cross. The circular path is associated with rotations of the string y . If the starting point is known, the only occurrence of letter \mathbf{b} here, following the path produces the initial string y	39
1.27	Sizes of texts compressed with three algorithms.	39

Chapter 1

Text data compression algorithms

MAXIME CROCHEMORE, Gaspard Monge Institute, Université de Marne-la-Vallée, France
and King's College London, UK

THIERRY LECROQ, Laboratoire d'Informatique, de Traitement de l'Information et des
Systèmes, Université de Rouen, France

1.1 Text compression

The chapter describes a few algorithms that compress texts. Compression serves both to save storage space and to save transmission time. We shall assume that the text is stored in a file. The aim of compression algorithms is to produce a new file, as short as possible, containing the compressed version of the same text. Methods presented here reduce the representation of text without any loss of information, so that decoding the compressed text restores exactly the original data.

The term “text” should be understood in a wide sense. It is clear that texts can be written in natural languages or can be texts usually generated by translators (like various types of compilers). But texts can also be images or other kinds of structures as well provided the data are stored in linear files. Texts considered here are sequence of characters from an finite alphabet Σ of size σ .

The interest in data compression techniques remains important even if mass storage systems improve regularly because the amount of data grows accordingly. Moreover, a con-

sequence of the extension of computer networks is that the quantity of data they exchange grows exponentially. So, it is often necessary to reduce the size of files to reduce proportionally their transmission times. Other advantages in compressing files regard two connected issues: integrity of data and security. While the first is easily accomplished through redundancy checks during the decompression phase, the second often requires data compression before applying cryptography.

This chapter contains three classical text compression algorithms. Variants of these algorithms are implemented in practical compression software, in which they are often combined together or with other elementary methods. Moreover, we present all-purpose methods, that is, methods in which no sophisticated modeling of the statistics of texts is done. An adequate modeling of a well-defined family of texts may increase significantly the compression when coupled with the coding algorithms of this chapter.

Compression ratios of the methods depend on the input data. However, most often, the size of compressed text vary from 30% to 50% of the size of the input. At the end of this chapter, we present ratios obtained by the methods on several example texts. Results of this type can be used to compare the efficiency of methods in compressing data. But, the efficiency of algorithms is also evaluated by their running times, and sometimes by the amount of memory space they require at run time. These elements are important criteria of choice when a compression algorithm is to be implemented in a telecommunication software.

Two strategies are applied to design the algorithms. The first strategy is a statistical method that takes into account the frequencies of symbols to build a uniquely decipherable code optimal with respect to the compression (Sections 1.2 and 1.3). Section 1.4 presents a refinement of the coding algorithm of Huffman based on the binary representation of numbers. Huffman codes contain new codewords for the symbols occurring in the text. In this method fixed-length blocks of bits are encoded by different codewords. *A contrario* the second strategy encodes variable-length segments of the text (Section 1.5). To put it simply, the algorithm,

while scanning the text, replaces some already read segments by just a pointer to their first occurrences. This second strategy often provides better compression ratios.

1.2 Static Huffman coding

The Huffman method is an optimal statistical coding. It transforms the original code used for characters of the text (ASCII code on 8 bits, for instance). Coding the text is just replacing each symbol (more exactly each occurrence of it) by its new **codeword**. The method works for any length of blocks (not only 8 bits), but the running time grows exponentially with the length. In the following, we assume that symbols are originally encoded on 8 bits to simplify the description.

The Huffman algorithm uses the notion of **prefix code**. A prefix code is a set of words containing no word that is a **prefix** of another word of the set. The advantage of such a code is that decoding is immediate. Moreover, it can be proved that this type of code does not weaken the compression.

A prefix code on the binary alphabet $\{0, 1\}$ corresponds to a binary tree in which the links from a node to its left and right children are labeled by **0** and **1** respectively. Such a tree is called a (digital) **trie**. Leaves of the trie are labeled by the original characters, and labels of branches are the words of the code (codewords of characters). Working with prefix code implies that codewords are identified with leaves only. Moreover, in the present method codes are complete: they correspond to complete tries, *i.e.* tree in which internal nodes have all exactly two children.

In the model where characters of the text are given new codewords, the Huffman algorithm builds a code that is optimal in the sense that the compression is the best possible (if the model of the source text is a zero-order Markov process, that is if the probability of symbol occurrence are independent). The length of the encoded text is minimum. The code depends on the input text, and more precisely on the frequencies of characters in the text. The most

```

H-COUNT (fin)
1  for each character  $a \in \Sigma$ 
2    do  $freq(a) \leftarrow 0$ 
3  while not end of file fin and  $a$  is the next symbol
4    do  $freq(a) \leftarrow freq(a) + 1$ 
5   $freq(\text{END}) \leftarrow 1$ 

```

Figure 1.1: Counts the character frequencies.

frequent characters are given shortest codewords while the least frequent symbols correspond to the longest codewords.

1.2.1 Encoding

The complete compression algorithm is composed of three steps: count of character frequencies, construction of the prefix code, encoding of the text. The last two steps use information computed by their preceding step.

The first step consists in counting the number of occurrences of each character in the original text (see Figure 1.1). We use a special end marker, denoted by **END**, which virtually appears only once at the end of the text. It is possible to skip this first step if fixed statistics on the alphabet are used. In this case however the method is optimal according to the statistics, but not necessarily for the specific text.

The second step of the algorithm builds the tree of a prefix code, called a Huffman tree, using the character frequency $freq(a)$ of each character a in the following way:

- create a one-node tree t for each character a , setting $weight(t) = freq(a)$ and $label(t) = a$,
- repeat
 - Extract the two least weighted trees t_1 and t_2 ,
 - Create a new tree t_3 having left subtree t_1 , right subtree t_2 , and weight $weight(t_3) = weight(t_1) + weight(t_2)$,

```

H-BUILD-TREE
1  for each  $a \in \Sigma \cup \{\text{END}\}$ 
2      do if  $\text{freq}(a) \neq 0$ 
3          then create a new node  $t$ 
4               $\text{weight}(t) \leftarrow \text{freq}(a)$ 
5               $\text{label}(t) \leftarrow a$ 
6   $l\text{leaves} \leftarrow$  list of all created nodes in increasing order of weight
7   $l\text{trees} \leftarrow$  empty list
8  while  $\text{LENGTH}(l\text{leaves}) + \text{LENGTH}(l\text{trees}) > 1$ 
9      do  $(\ell, r) \leftarrow$  extract two nodes of smallest weight (among the two nodes at the
          beginning of  $l\text{leaves}$  and the two nodes at the beginning of  $l\text{trees}$ )
10     create a new node  $t$ 
11      $\text{weight}(t) \leftarrow \text{weight}(\ell) + \text{weight}(r)$ 
12      $\text{left}(t) \leftarrow \ell$ 
13      $\text{right}(t) \leftarrow r$ 
14     insert  $t$  at the end of  $l\text{trees}$ 
16 return  $t$ 

```

Figure 1.2: Builds the Huffman coding tree.

- until only one tree remains.

The tree is constructed by the algorithm H-BUILD-TREE in Figure 1.2. The implementation uses two linear lists. The first list, $l\text{leaves}$, contains the leaves of the future tree associated each with a symbol. The list is sorted in increasing order of weights of leaves (frequencies of symbols). The second list, $l\text{trees}$, contains the newly created trees. The operation of extracting the two least weighted trees is done by checking the two first trees of the list $l\text{leaves}$ and the two first trees of the list $l\text{trees}$. Each new tree is inserted at the end of the list of the trees. The only tree remaining at the end of the procedure is the coding tree.

After the coding tree is built, it is possible to recover the codewords associated with characters by a simple depth-first-search of the tree (see Figure 1.3); $\text{codeword}(a)$ denotes the binary codeword associated with the character a .

In the third step, the original text is encoded. Since the code depends on the original text, in order to be able to decode the compressed text, the coding tree and the original

```

H-BUILD-CODE ( $t$ ,  $length$ )
1  if  $t$  is not a leaf
2      then  $temp[length] \leftarrow 0$ 
3          H-BUILD-CODE ( $left(t)$ ,  $length + 1$ )
4           $temp[length] \leftarrow 1$ 
5          H-BUILD-CODE ( $right(t)$ ,  $length + 1$ )
6  else  $codeword(label(t)) \leftarrow temp[0..length - 1]$ 

```

Figure 1.3: Builds character codewords from the coding tree.

```

H-ENCODE-TREE ( $fout$ ,  $t$ )
1  if  $t$  is not a leaf
2      then write a 0 in the file  $fout$ 
3          H-ENCODE-TREE ( $fout$ ,  $left(t)$ )
4          H-ENCODE-TREE ( $fout$ ,  $right(t)$ )
5  else write a 1 in the file  $fout$ 
6      write the original code of  $label(t)$  in the file  $fout$ 

```

Figure 1.4: Stores the coding tree in the compressed file.

codewords of symbols must be stored together with the compressed text.

This information is placed in a header of the compressed file, to be read at decoding time just before the decompression starts. The header is written during a depth-first traversal of the tree. Each time an internal node is encountered a **0** is produced. When a leaf is encountered a **1** is produced followed by the original code of the corresponding character on 9 bits (so that the end marker can be equal to 256 if all the 8-bit characters appear in the original text). This part of the encoding algorithm is shown in Figure 1.4.

After the header of the compressed file is made, the encoding of the original text is realized by the algorithm of Figure 1.5.

A complete implementation of the Huffman algorithm, composed of the three steps described above, is given in Figure 1.6.

Example 1.1:

```
H-ENCODE-TEXT (fin, fout)
1  while not end of file fin and a is the next symbol
2      do write codeword(a) in the file fout
3  write codeword(END) in the file fout
```

Figure 1.5: Encodes the characters in the compressed file.

```
H-ENCODING (fin, fout)
1  H-COUNT (fin)
2   $t \leftarrow$  H-BUILD-TREE
3  H-BUILD-CODE ( $t$ , 0)
4  H-ENCODE-TREE (fout,  $t$ )
5  H-ENCODE-TEXT (fin, fout)
```

Figure 1.6: Complete function for Huffman encoding.

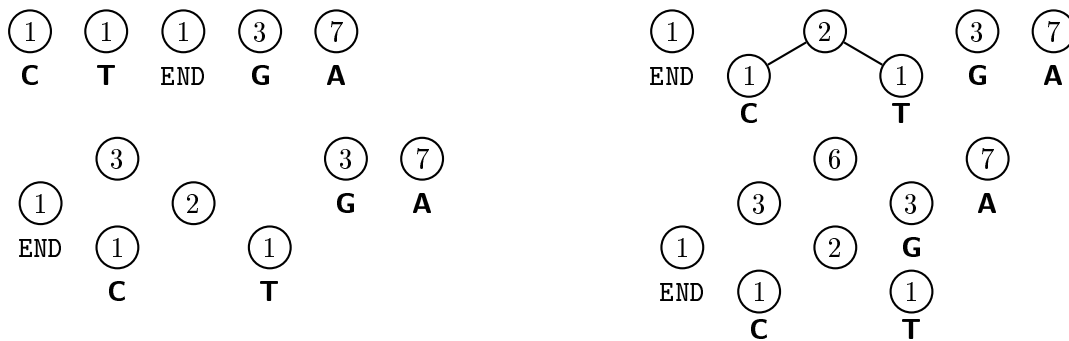
$y = \mathbf{ACAGAATAGAGA}$

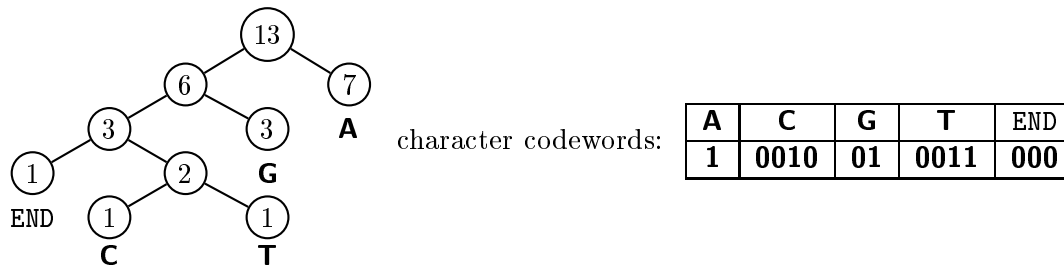
Length of $y = 12 \times 8 = 104$ bits (assuming an 8-bit code)

Character frequencies:

A	C	G	T	END
7	1	3	1	1

Different steps during the construction of the coding tree:





Encoded tree: **0001***binary*(END,9)**01***binary*(C,9)**1***binary*(T,9)**1***binary*(G,9)**1***binary*(A,9),

which produces a header of length 54 bits:

0001 10000000 01 001000011 1 001010100 1 001000111 1 001000001

Encoded text:

$\underbrace{1}_{\mathbf{A}}$
 $\underbrace{0010}_{\mathbf{C}}$
 $\underbrace{1}_{\mathbf{A}}$
 $\underbrace{01}_{\mathbf{G}}$
 $\underbrace{1}_{\mathbf{A}}$
 $\underbrace{1}_{\mathbf{A}}$
 $\underbrace{0011}_{\mathbf{T}}$
 $\underbrace{1}_{\mathbf{A}}$
 $\underbrace{01}_{\mathbf{G}}$
 $\underbrace{1}_{\mathbf{A}}$
 $\underbrace{01}_{\mathbf{G}}$
 $\underbrace{1}_{\mathbf{A}}$
 $\underbrace{000}_{\mathbf{END}}$

of length 24 bits

Total length of the compressed file: 78 bits. □

The construction of the tree takes $O(\sigma \log \sigma)$ time if the sorting of the list of the leaves is implemented efficiently. The rest of the encoding process runs in time linear in the sum of the sizes of the original and compressed texts.

1.2.2 Decoding

Decoding a file containing a text compressed by Huffman algorithm is a mere programming exercise. First, the coding tree is rebuilt by the algorithm of Figure 1.7. Then, the original text is recovered by parsing the compressed text with the coding tree. The process begins at the root of the coding tree, and follows a left edge when a **0** is read or a right edge when a **1** is read. When a leaf is encountered, the corresponding character (in fact the original codeword of it) is produced and the parsing resumes at the root of the tree. The process ends when the codeword of the end marker is encountered. An implementation of the decoding of the text is

```

H-REBUILD-TREE (fin, t)
1  read bit b from fin
2  if b = 1
3      then make t a leaf
4          label(t) ← symbol corresponding to the 9 next bits read from fin
5  else create a new node ℓ
6      left(t) ← ℓ
7      H-REBUILD-TREE (fin, ℓ)
8      create a new node r
9      right(t) ← r
10     H-REBUILD-TREE (fin, r)

```

Figure 1.7: Rebuilds the tree from the header of compressed file.

```

H-DECODE-TEXT (fin, fout, root)
1  t ← root
2  while label(t) ≠ END
3      do if t is a leaf
4          then write label(t) in the file fout
5              t ← root
6      else read bit b from fin
7          if b = 1
8              then t ← right(t)
9              else t ← left(t)

```

Figure 1.8: Recovers the original text.

```

H-DECODING ( $fin, fout$ )
1  create a new node  $root$ 
2  H-REBUILD-TREE ( $fin, root$ )
3  H-DECODE-TEXT ( $fin, fout, root$ )

```

Figure 1.9: Complete function for Huffman decoding.

presented in Figure 1.8.

The complete decoding program is given in Figure 1.9. It calls the preceding functions. The running time of the decoding program is linear in the sum of the sizes of the texts it manipulates.

1.3 Dynamic Huffman coding

The two main drawbacks of the static Huffman method are: first, if the frequencies of characters the source text are not known *a priori*, the source text has to be read twice; second, the coding tree must be included in the compressed file. This is avoided by a dynamic method where the coding tree is updated each time a symbol is read from the source text. The current tree is a Huffman tree related to the part of the text that is already treated. The tree evolves exactly in the same way during the decoding process. The efficiency of the method is based on a characterization of Huffman trees, known as the *siblings property*.

Siblings property: Let T be a Huffman tree with n leaves (a complete binary weighted tree built by the procedure H-BUILD-TREE in which all leaves have positive weights). Then the nodes of T can be arranged in a sequence $(x_0, x_1, \dots, x_{2n-2})$ such that:

1. the sequence of weights $(weight(x_0), weight(x_1), \dots, weight(x_{2n-2}))$ is in decreasing order;
2. for any i ($0 \leq i \leq n - 2$), the consecutive nodes x_{2i+1} and x_{2i+2} are siblings (they have the same parent).

The compression and decompression processes initialize the dynamic Huffman tree by a one-node tree that correspond to an artificial character, denoted by **ART**. The weight of this single node is 1.

1.3.1 Encoding

Each time a symbol a is read from the source text, its codeword in the tree is sent. However this happens only if a appeared previously. Otherwise the code of **ART** is sent followed by the original codeword of a . Afterwards, the tree is modified in the following way: first, if a never occurred before, a new internal node is created and its two children are a new leaf labeled by a and the leaf **ART**; then, the tree is updated (see below) to get a Huffman tree for the new prefix of text.

Implementation

Each node is identified with an integer n , the root is the integer 0. The invariant of compression and decompression algorithms is that, if the tree has m nodes, the sequence of nodes $(m - 1, \dots, 1, 0)$ satisfies the siblings property. The tree is stored in a table, and we use the next notations, for a node n :

- $parent(n)$ is the parent of n ($parent(root) = \text{UNDEFINED}$),
- $child(n)$ is the left child of n (if n is an internal node, otherwise $child(n) = \text{UNDEFINED}$), and $child(n) + 1$ is its right child (this is useful only at decoding time),
- $label(n)$ is the symbol associated with n when n is a leaf,
- $weight(n)$ is the weight of n (it is the frequency of $label(n)$ if n is a leaf).

Indeed, the child link is useful only at decoding time so that the implementation may differ between the two phases of the algorithm. But, to simplify the description and give a uniform treatment of the data structure, we assume the same implementation during the encoding and the decoding steps. Weights of nodes are handled by the procedure **DH-UPDATE**.

```

DH-INIT
1  root ← 0
2  child(root) ← UNDEFINED
3  weight(root) ← 1
4  parent(root) ← UNDEFINED
5  for each letter a ∈ Σ ∪ {END}
6      do leaf[a] ← UNDEFINED
7  leaf[ART] ← root

```

Figure 1.10: Initializes the dynamic Huffman tree.

```

DH-ENCODING (fin, fout)
1  DH-INIT
2  while not end of file fin and a is the next symbol
3      do DH-ENCODE-SYMBOL (a, fout)
4          DH-UPDATE (a)
5  DH-ENCODE-SYMBOL (END, fout)

```

Figure 1.11: Complete function for dynamic Huffman encoding.

The correspondence between symbols and leaves of the tree is done by a table called *leaf*: for each symbol $a \in \Sigma \cup \{\text{END}\}$, *leaf*[*a*] is the corresponding leaf of the tree, if any.

The coding tree initially contains only one node labeled by symbol ART. The initialization is given in Figure 1.10.

Encoding the source text is a succession of three steps: read a symbol *a* from the source text, encode the symbol *a* according to the current tree, update the tree. It is described in Figure 1.11.

Encoding a symbol *a* already encountered consists in accessing its associated leaf *leaf*[*a*], then computing its codeword by a bottom-up walk in the tree (following the *parent* links up to the root). Each time a node *n* (\neq root) is encountered if *n* is odd a 1 is sent (*n* is the right child of its parent), and if *n* is even a 0 is sent (*n* is then the left child of its parent). As the codeword of *a* is read in reverse direction a stack *S* is used to temporarily store the bits and

```

DH-ENCODE-SYMBOL (a, fout)
1  S ← empty stack
2  n ← leaf[a]
3  if n = UNDEFINED
4      then n ← leaf[ART]
5  while n ≠ root
6      do if n is odd
7          then PUSH (S, 1)
8          else PUSH (S, 0)
9          n ← parent(n)
10 SEND (S, fout)
11 if leaf[a] = UNDEFINED
12     then write in fout the original codeword of a on 9 bits
13     DH-ADD-NODE (a)

```

Figure 1.12: Encodes one symbol.

```

DH-ADD-NODE (a)
1  transform leaf[ART] into
2     an internal node of weight 1 with
3     a left child of weight 0 for leaf[a],
4     a right child of weight 1 for leaf[ART].

```

Figure 1.13: Adds a new symbol in the tree.

```

DH-DECODING (fin, fout)
1  DH-INIT
2   $a \leftarrow$  DH-DECODE-SYMBOL (fin)
3  while  $a \neq$  END
4      do write  $a$  in fout
5          DH-UPDATE ( $a$ )
6           $a \leftarrow$  DH-DECODE-SYMBOL (fin)

```

Figure 1.14: Complete function for dynamic Huffman decoding.

send them properly. The procedure SEND ($S, fout$) send the bits of the stack S in the correct order to the file *fout*.

If a has not been encountered yet, then the code of **ART** is sent followed by the original codeword of a on 9 bits, and a new leaf is created for a (see Figures 1.12 and 1.13).

1.3.2 Decoding

At decoding time the compressed text is parsed with the coding tree. The current node is initialized with the root like in the encoding algorithm, and then the tree evolves symmetrically. Each time a **0** is read from the compressed file the walk down the tree follows the left link, and it follows the right link if a **1** is read. When the current node is a leaf, its associated symbol is written in the output file and the tree is updated exactly as it is done during the encoding phase.

Implementation

As for the encoding process the current Huffman tree is stored in a table, and it is initialized with the artificial symbol **ART**. The same elements of the data structure are used during the decompression. Note that the next node when parsing a bit b from node n is just $child(n) + b$ with the convention on left-right links and 0-1 bits. The tree is updated by the procedure DH-UPDATE used previously and that is described in the next section. Figures 1.14 and 1.15 display the decoding mechanism.

DH-DECODE-SYMBOL (fin)

```

1   $n \leftarrow root$ 
2  while  $child(n) \neq UNDEFINED$ 
3      do read bit  $b$  from  $fin$ 
4           $n \leftarrow child(n) + b$ 
5   $a \leftarrow label(n)$ 
6  if  $a = ART$ 
7      then  $a \leftarrow$  symbol corresponding to the next 9 bits read from  $fin$ 
8          DH-ADD-NODE ( $a$ )
9  return ( $a$ )

```

Figure 1.15: Decodes one symbol from the input file.

DH-UPDATE (a)

```

1   $n \leftarrow leaf[a]$ 
2  while  $n \neq root$ 
3      do  $weight(n) \leftarrow weight(n) + 1$ 
4           $m \leftarrow n$ 
5          while  $weight(m - 1) < weight(n)$ 
6              do  $m \leftarrow m - 1$ 
7          DH-SWAP-NODES ( $m, n$ )
8           $n \leftarrow parent(m)$ 
9   $weight(root) \leftarrow weight(root) + 1$ 

```

Figure 1.16: Updates the current Huffman tree.

1.3.3 Updating

During encoding and decoding phases the current tree has to be updated to take into account the correct frequency of symbols. When a new symbol is considered the weight of its associated leaf is incremented by 1, and the weights of ancestors have to be modified correspondingly. The procedure that realizes the operation is shown in Figures 1.16. Its proof of correctness is based on the siblings property.

We explain how the procedure DH-UPDATE works. First, the weight of the leaf n corresponding to a is incremented by 1. Then, if point 1 of the siblings property is no longer satisfied, node n is exchanged with the closest node m ($m < n$) in the list such that $weight(m) <$

$weight(n)$. Doing so, the nodes remain in decreasing order of their weights. It is important here that leaves have positive weights because this guarantees that m is not a parent nor an ancestor of node n . Afterwards, the same operation is repeated on the parent of n until the root of the tree is reached.

The aim of procedure DH-SWAP-NODES used in Figures 1.16 is to exchange the subtrees rooted at its input nodes m and n . In concrete terms, this remains to exchange the records stored at the two nodes in the table. It is meant that nothing is to be done if $m = n$.

Example 1.2:

$y = \mathbf{ACAGAATAGAGA}$

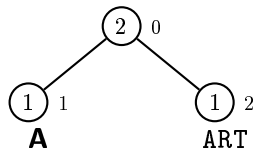
Initial tree:



Next symbol is **A**:

The ASCII code of **A** is sent on 9 bits;

Bits sent: **001000001**

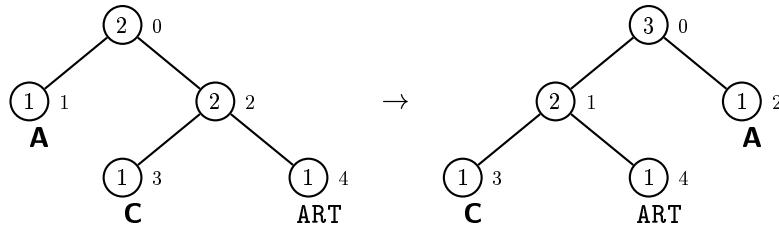


Next symbol is **C**:

The code of ART is sent followed by the ASCII code of **C** on 9 bits;

Bits sent: **1 001000011**

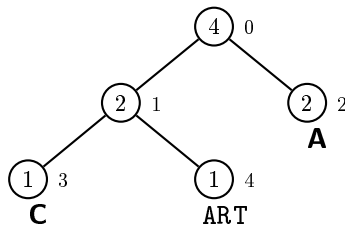
Nodes 1 and 2 are swapped.



Next symbol is **A**:

The code of **A** is sent;

Bit sent: **1**

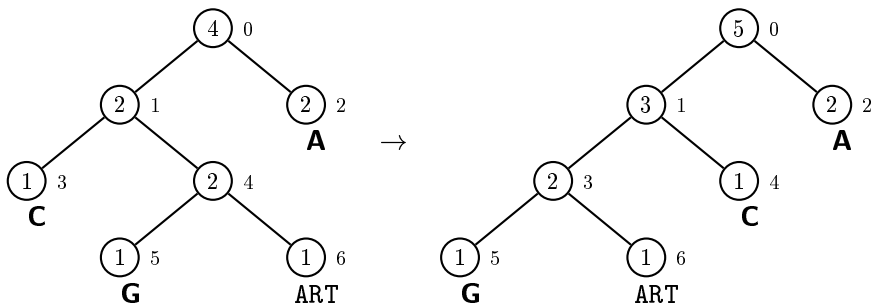


Next symbol is **G**:

The code of **ART** is sent followed by the ASCII code of **G** on 9 bits;

Bits sent: **01 001000111**

Nodes 3 and 4 are swapped.



Finally, the entire sequence of bits sent is the following:

$\underbrace{001000001}_{\mathbf{A}}$
 $\underbrace{1}_{\mathbf{ART}}$
 $\underbrace{001000011}_{\mathbf{C}}$
 $\underbrace{1}_{\mathbf{A}}$
 $\underbrace{01}_{\mathbf{ART}}$
 $\underbrace{001000111}_{\mathbf{G}}$
 $\underbrace{1}_{\mathbf{A}}$
 $\underbrace{1}_{\mathbf{A}}$
 $\underbrace{101}_{\mathbf{ART}}$

$\underbrace{001010100}_{\mathbf{T}}$
 $\underbrace{0}_{\mathbf{A}}$
 $\underbrace{100}_{\mathbf{G}}$
 $\underbrace{0}_{\mathbf{A}}$
 $\underbrace{100}_{\mathbf{G}}$
 $\underbrace{0}_{\mathbf{A}}$
 $\underbrace{111}_{\mathbf{ART}}$
 $\underbrace{100000000}_{\mathbf{END}}$

The total length of the compressed text is 66.

□

```

AR-ENCODE (fn)
1   $\ell \leftarrow 0$ 
2   $h \leftarrow 1$ 
3  while not end of file fn and  $a_i$  is the next symbol
4      do   $\ell \leftarrow \ell + (h - \ell) * \ell_i$ 
5           $h \leftarrow \ell + (h - \ell) * h_i$ 
6  return( $\ell$ )

```

Figure 1.17: Basic arithmetic encoding.

1.4 Arithmetic coding

1.4.1 Encoding

The basic idea of arithmetic coding is to consider symbol as digits of a numeration system, and texts as decimal parts of numbers between 0 and 1. The length of the interval attributed to a digit (it is 0.1 for digits in the usual base 10 system) is made proportional to the frequency of the digit in the text. The encoding is thus assimilated to a change in the base of a numeration system. To cope with precision problems, the number corresponding to a text is handled via a lower bound and an upper bound, which remains to associate with a text a sub-interval of $[0, 1[$. The compression comes from the fact that large intervals require less precision to separate their bounds.

More formally, the interval associated with each symbol $a_i \in \Sigma$ ($1 \leq i \leq \sigma$) is denoted $I(a_i) = [\ell_i, h_i[$. The intervals satisfy the conditions: $\ell_1 = 0$, $h_\sigma = 1$, and $\ell_i = h_{i-1}$ for $1 < i \leq \sigma$. Note that $I(a_i) \cap I(a_j) = \emptyset$ if $a_i \neq a_j$.

The encoding phase consists in computing the interval corresponding to the input text. The basic step that deals with a symbol a_i of the source text transforms the current interval $[\ell, h[$ into $[\ell', h'[$ where $\ell' = \ell + (h - \ell) * \ell_i$ and $h' = \ell + (h - \ell) * h_i$, starting with the initial interval $[0, 1[$ (see Figure 1.17). Indeed, in a theoretical approach, ℓ only is needed to encode the input text.

```

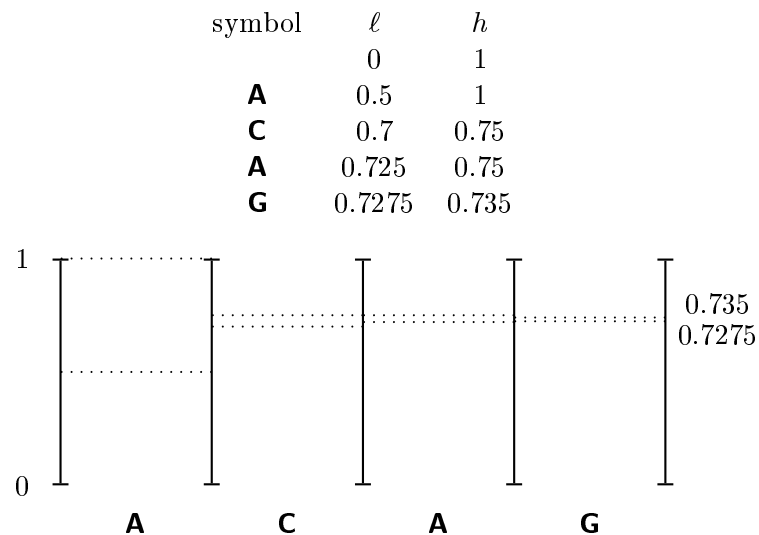
AR-DECODE ( $\ell, fout$ )
1  while  $\ell \neq 0$ 
2      do  find  $a_i$  such that  $\ell \in I(a_i)$ 
3          write  $a_i$  in file  $fout$ 
4           $\ell \leftarrow (\ell - \ell_i)/(h_i - \ell_i)$ 

```

Figure 1.18: Basic arithmetic decoding.

Example 1.3:
 $\Sigma = \{\mathbf{A}, \mathbf{C}, \mathbf{G}, \mathbf{T}\}, \sigma = 4$

$$I(\mathbf{A}) = [0.5, 1[, \quad I(\mathbf{C}) = [0.4, 0.5[, \quad I(\mathbf{G}) = [0.1, 0.4[, \quad I(\mathbf{T}) = [0, 0.1[$$

Encoding **ACAG** gives:

□

1.4.2 Decoding

The reverse operation of decoding a text compressed by the previous algorithm theoretically requires only the lower bound ℓ of the interval. Decoding the number ℓ is done as follows: first find the symbol a_i such that $\ell \in I(a_i)$, produced the symbol a_i , and then replace ℓ by ℓ' defined

by:

$$\ell' \leftarrow \frac{\ell - \ell_i}{h_i - \ell_i}.$$

The same process is repeated until $\ell = 0$ (see Figure 1.18). Indeed, the implementation of the decoding phase simulates what is done on the current interval considered at encoding time.

Example 1.4:

$$\Sigma = \{\mathbf{A}, \mathbf{C}, \mathbf{G}, \mathbf{T}\}, \sigma = 4$$

$$I(\mathbf{A}) = [0.5, 1[, \quad I(\mathbf{C}) = [0.4, 0.5[, \quad I(\mathbf{G}) = [0.1, 0.4[, \quad I(\mathbf{T}) = [0, 0.1[$$

Decoding $\ell = 0.7275$:

ℓ	a_i
0.7275	A
0.455	C
0.55	A
0.1	G
0	

□

1.4.3 Implementation

The main problem when implementing the arithmetic coding compression algorithm is to cope with precision on real numbers operations. The $[0, 1[$ interval of real numbers is substituted by the interval of integers $[0, 2^N - 1[$, where N is a fixed integer.

So, the algorithms work with integral values of size N . During the process, each time the binary representation of bounds ℓ and h have a common prefix this prefix is sent and ℓ is shifted to the left and filled with 0's while h is shifted to the left and filled with 1's.

The intervals associated with symbols of the alphabet are computed with the help of symbol frequencies, in a dynamic way: each character frequency is initialized with 1 and is incremented each time the symbol is encountered.

```

AR-ENCODING (fin, fout)
1  Initialize tables freq and cum-freq
2   $\ell \leftarrow 0$ 
3   $h \leftarrow 2^N - 1$ 
4  waiting-counter  $\leftarrow 0$ 
5  while not end of file fin and  $a_i$  is the next symbol
6      do AR-ENCODE-SYMBOL ( $a_i$ )
7          Update tables freq and cum-freq
8          Maintain symbols in decreasing order of frequencies
9  AR-ENCODE-SYMBOL (END)
10 waiting-counter  $\leftarrow$  waiting-counter + 1
11 AR-SEND-BIT (leftmost bit of  $h$ )

```

Figure 1.19: Complete arithmetic encoding function.

```

AR-ENCODE-SYMBOL ( $a_i$ )
1   $\ell \leftarrow \ell + ((h - \ell + 1) * \text{cum-freq}[i]) / \text{cum-freq}[0]$ 
2   $h \leftarrow \ell + ((h - \ell + 1) * \text{cum-freq}[i - 1]) / \text{cum-freq}[0] - 1$ 
3  repeat
4      if leftmost bit of  $\ell$  = leftmost bit of  $h$ 
5          then AR-SEND-BIT (leftmost bit of  $h$ )
6               $\ell \leftarrow 2 * \ell$ 
7               $h \leftarrow 2 * h + 1$ 
8          else if  $h - \ell < \text{cum-freq}[0]$ 
9              then  $\ell \leftarrow 2 * (\ell - 2^{N-2})$ 
10                  $h \leftarrow 2 * (h - 2^{N-2}) + 1$ 
11                 waiting-counter  $\leftarrow$  waiting-counter + 1
12 until leftmost bit of  $\ell \neq$  leftmost bit of  $h$  and  $h - \ell \geq \text{cum-freq}[0]$ 

```

Figure 1.20: Encodes one symbol.

```

AR-SEND-BIT (bit, fout)
1  write bit in fout
2  while waiting-counter > 0
3      do if bit = 0
4          then write 1 in fout
5              write 0 in fout
6          waiting-counter  $\leftarrow$  waiting-counter - 1

```

Figure 1.21: Sends one bit followed by the waiting bits, if any.

We denote by $freq[i]$ the frequency of symbol a_i of the alphabet. We also consider the cumulative frequency of symbols, and set $cum-freq[i] = \sum_{j=i+1}^{\sigma} freq[j]$. Then $cum-freq[0]$ is the cumulative frequency of all symbols. Note that $cum-freq[0] - \sigma - 1$ is the length of the prefix of the input scanned so far. The symbols are maintained in decreasing order of their frequencies. This obviously save on the expected number of operations to update the table $cum-freq$.

Then when a symbol a_i is read from the source text, the current interval $[\ell, h[$ is updated in the following way:

$$\ell \leftarrow \ell + \frac{(h - \ell + 1) * cum-freq[i]}{cum-freq[0]}$$

$$h \leftarrow \ell + \frac{(h - \ell + 1) * cum-freq[i - 1]}{cum-freq[0]} - 1$$

The common prefix (if any) of ℓ and h is sent and ℓ and h are shifted and filled (respectively with 0's and 1's). At this point, if the interval is too short (if $h - \ell < cum-freq[0]$) it is extended to $[2 * (\ell - 2^{N-2}), 2 * (h - 2^{N-2}) + 1[$ and a waiting counter is incremented by 1. And the same operation is repeated as long as the interval is too short. After that, when a bit is sent, the reverse bit is sent the number of times indicated by the waiting counter.

Example 1.5:

$$\Sigma = \{\mathbf{A}, \mathbf{C}, \mathbf{G}, \mathbf{T}\}, N = 8$$

		A	C	G	T	END
i	0	1	2	3	4	5
$cum-freq$	5	4	3	2	1	0
$freq$	0	1	1	1	1	1

$$[0, 255[\xrightarrow{\mathbf{A}} [0 + \frac{(255-0+1)*4}{5}, 0 + \frac{(255-0+1)*5}{5} - 1[= [204, 255[= [11001100_2, 11111111_2[$$

11 is sent and $[00110000_2, 11111111_2[= [48, 255[$ is the next interval;

		A	C	G	T	END
i	0	1	2	3	4	5
$cum-freq$	6	4	3	2	1	0
$freq$	0	2	1	1	1	1

$$[48, 255[\xrightarrow{\mathbf{C}} [96, 231[= [10011000_2, 10111001_2[$$

10 is sent and $[01100000_2, 11100111_2[= [96, 231[$ is the next interval;

		A	C	G	T	END
<i>i</i>	0	1	2	3	4	5
<i>cum-freq</i>	7	5	3	2	1	0
<i>freq</i>	0	2	2	1	1	1

$[96, 231[\xrightarrow{\mathbf{A}}$ $[193, 231[= [11000001_2, 11100111_2[$

11 is sent and $[00000100_2, 10011111_2[= [4, 159[$ is the next interval;

Next symbol is **G**: **001** is sent;

Next symbol is **A**: **1** is sent;

Next symbol is **A**: **1** is sent;

Next symbol is **T**: **0111** is sent;

Next symbol is **A**: **10** is sent;

Next symbol is **G**: nothing is sent and the current interval becomes $[111, 139[$;

$[111, 139[\xrightarrow{\mathbf{A}}$ $[127, 139[= [01111111_2, 10001011_2[$, nothing is sent

The interval is too short and replaced by $[2 * (127 - 2^{N-2}), 2 * (139 - 2^{N-2}) + 1[= [126, 151[$

and one bit is waiting, $[01111110, 10010111[= [126, 151[$ is the next interval;

$[126, 151[\xrightarrow{\mathbf{G}}$ $[134, 138[= [10000110_2, 10001010_2[$

1+0+000 are sent and $[01100000_2, 10101111_2[= [96, 175[$ is the next interval;

$[96, 175[\xrightarrow{\mathbf{A}}$ $[141, 175[= [10001101_2, 10101111_2[$

10 is sent and $[00110100_2, 10111111_2[= [52, 191[$ is the next interval;

$[52, 191[\xrightarrow{\text{END}}$ $[52, 59[= [00110100_2, 00111011_2[$

0011 is sent;

```

AR-DECODING (fin, fout)
1  Init the tables freq and cum-freq
2  value ← 0
3  for i ← 1 to N
4      do read bit b from fin
5          value ← 2 * value + b
6  ℓ ← 0
7  h ← 2N − 1
8  repeat
9      i ← AR-DECODE-SYMBOL (fin)
10     if ai ≠ END
11         then write ai in fout
12             Update the tables freq and cum-freq
13             Maintain symbols in decreasing order of frequencies
14 until ai = END

```

Figure 1.22: Complete arithmetic decoding function.

1+0 are sent in order to finish the encoding process. □

The decoding process is exactly the reverse of the coding process. It uses a window of size N on the compressed file. First, the window is filled with the first N bits of the compressed file and $value$ is the corresponding base 2 number. The current interval is initialized with $\ell = 0$ and $h = 2^N - 1$.

Then, the symbol a_i to be produced is the first character such that:

$$cum\text{-}freq[i] > \frac{(value - \ell + 1) * cum\text{-}freq[0] - 1}{h - \ell + 1},$$

and ℓ and h are then updated exactly in the same way than during the coding process. If the binary representations of ℓ and h have a common prefix of length p they are both shifted p binary position to the left and ℓ is filled by 0's, h is filled with 1's. The window on the compressed file is shifted p symbols to the right and the variable $value$ is updated correspondingly. The tables $freq$ and $cum\text{-}freq$ are updated and the symbols are maintained in decreasing order of the frequencies as in the coding process.

```

AR-DECODE-SYMBOL (fin)
1   $cum \leftarrow ((value - \ell + 1) * cum-freq[0] - 1) / (h - \ell + 1)$ 
2   $i \leftarrow 1$ 
3  while  $cum-freq[i] > cum$ 
4      do  $i \leftarrow i + 1$ 
5   $\ell \leftarrow \ell + ((h - \ell + 1) * cum-freq[i]) / cum-freq[0]$ 
6   $h \leftarrow \ell + ((h - \ell + 1) * cum-freq[i - 1]) / cum-freq[0] - 1$ 
7  repeat
8      if leftmost bit of  $\ell =$  leftmost bit of  $h$ 
9          then  $\ell \leftarrow 2 * \ell$ 
10              $h \leftarrow 2 * h + 1$ 
11             read bit  $b$  from fin
12              $value \leftarrow 2 * value + b$ 
13      else if  $h - \ell < cum-freq[0]$ 
14          then  $\ell \leftarrow 2 * (\ell - 2^{N-2})$ 
15              $h \leftarrow 2 * (h - 2^{N-2}) + 1$ 
16             read bit  $b$  from fin
17              $value \leftarrow 2 * (value - 2^{N-2}) + b$ 
18  until leftmost bit of  $\ell \neq$  leftmost bit of  $h$  and  $h - \ell \geq cum-freq[0]$ 
19  return  $i$ 

```

Figure 1.23: Decodes one symbol.

This is repeated until the symbol END is produced.

Example 1.6:

Decoding the text **111011001110111101000010001110**

$\Sigma = \{\mathbf{A}, \mathbf{C}, \mathbf{G}, \mathbf{T}\}$

		A	C	G	T	END
i	0	1	2	3	4	5
$cum\text{-}freq$	5	4	3	2	1	0
$freq$	0	1	1	1	1	1

$value = 236 = 11101100_2$

$cum = 4$

$[0, 255[\xrightarrow{\mathbf{A}} [204, 255[= [11001100_2, 11111111_2[: \text{shift by 2}$

The next interval is $[00110000_2, 11111111_2[$ and $value = 10110011_2 = 179$

		A	C	G	T	END
i	0	1	2	3	4	5
$cum\text{-}freq$	6	4	3	2	1	0
$freq$	0	2	1	1	1	1

$value = 179 = 10110011_2$

$cum = 3$

$[48, 255[\xrightarrow{\mathbf{C}} [152, 185[= [10011000_2, 10111001_2[: \text{shift by 2}$

The next interval is $[01100000_2, 11100111_2[$ and $value = 11001110_2 = 206$

		A	C	G	T	END
i	0	1	2	3	4	5
$cum\text{-}freq$	7	5	3	2	1	0
$freq$	0	2	2	1	1	1

$value = 206 = 11001110_2$

$cum = 5$

$[96, 231[\xrightarrow{\mathbf{A}} [193, 231[= [11000001_2, 11100111_2[: \text{shift by 2}$

The next interval is $[00000100_2, 10011111_2[$ and $value = 00111011_2 = 59$

Next symbols are **GAATAG** and the current interval becomes $[111, 139[$ and $value = 132 = 10000100_2$

$cum = 10$

$[111, 139[\xrightarrow{\mathbf{A}}$ $[127, 139[= [01111111_2, 10001011_2[$: no shift

The interval is too short and is replaced by $[01111110_2, 10010111_2[$ and $value = 10001000_2 = 136$

$value = 136 = 10001000_2$

$cum = 6$

$[126, 151[\xrightarrow{\mathbf{G}}$ $[134, 138[= [10000110_2, 10001010_2[$: shift by 4

The next interval is $[01100000_2, 10101111_2[$ and $value = 10001110_2 = 142$

$value = 142 = 10001110_2$

$cum = 9$

$[96, 175[\xrightarrow{\mathbf{A}}$ $[141, 175[= [10001101_2, 10101111_2[$: shift by 2

The next interval is $[00110100_2, 10111111_2[$ and $value = 00111000_2 = 56$

$value = 56 = 00111000_2$

$cum = 0$

The symbol is **END**, the decoding process is over. □

Maintaining the symbols in decreasing order of frequencies can be done in $O(\log \sigma)$ using a suitable data structure (see Fenwick 1994).

1.5 LZW Coding

Ziv and Lempel designed a compression method using encoding **segments**. These segments of the original text are stored in a dictionary that is built during the compression process. When a segment of the dictionary is encountered later while scanning the text it is substituted by its index in the dictionary. In the model where portions of the text are replaced by pointers on previous occurrences, the Ziv-Lempel compression scheme can be proved to be asymptotically optimal (on large enough texts satisfying good conditions on the probability distribution of symbols).

The dictionary is the central point of the algorithm. It has the property of being prefix-closed (every prefix of a word of the dictionary is in the dictionary), so that it can be implemented efficiently as a trie. Furthermore, a hashing technique makes its implementation efficient. The version described in this section is called the Lempel-Ziv-Welsh method after several improvements introduced by Welsh. The algorithm is implemented by the `compress` command existing under the UNIX operating system.

1.5.1 Encoding

We describe the scheme of the coding method. The dictionary is initialized with all strings of length 1, the characters of the alphabet. The current situation is when we have just read a segment w of the text. Let a be the next symbol (just following the given occurrence w). Then we proceed as follows:

- If wa is not in the dictionary, we write the index of w in the output file, and add wa to the dictionary. We then reset w to a and process the next symbol (following a).
- If wa is in the dictionary we process the next symbol, with segment wa instead of w .

Initially, the segment w is set to the first symbol of the source text, so that it is clear that “ w belongs to the dictionary” is an invariant of the operations described above.

Example 1.7:

The alphabet is the 8-bit ASCII alphabet, $y = \mathbf{ACAGAATAGAGA}$

The dictionary initially contains the ASCII symbols, their indexes are their ASCII codewords.

A	C	A	G	A	A	T	A	G	A	G	A	w	written	added
	↑											A	65	AC , 257
		↑										C	67	CA , 258
			↑									A	65	AG , 259
				↑								G	71	GA , 260
					↑							A	65	AA , 261
						↑						A	65	AT , 262
							↑					T	84	TA , 263
								↑				A		
									↑			AG	259	AGA , 264
										↑		A		
											↑	AG		
											↑	AGA	264	
													256	

□

1.5.2 Decoding

The decoding method is symmetric to the coding algorithm. The dictionary is recovered while the decompression process runs. It is basically done in this way:

- read a code c in the compressed file,
- write in the output file the segment w having index c in the dictionary,
- add the word wa to the dictionary where a is the first letter of the next segment.

In this scheme, the dictionary is updated after the next segment is decoded because we need its first symbol a to concatenate at the end of the current segment w . So, a problem occurs if the next index computed at encoding time is precisely the index of the segment wa . Indeed, this happens only in a very special case, in which the symbol a is also the first symbol of w itself. This arises if the rest of the text to encode starts with a segment $azazax$ (a a symbol, z a string) for which az belongs to the dictionary but aza does not. During the compression process the index of az is output, and aza is added to the dictionary. Next, aza is read and

its index is output. During the decompression process the index of *aza* is read while the first occurrence of *az* has not been completed yet, the segment *aza* is not already in the dictionary. However, since this is the unique case where the situation occurs, the segment *aza* is recovered by taking the last segment *az* added to the dictionary concatenated with its first letter *a*.

Example 1.8:

Decoding the sequence 65, 67, 65, 71, 65, 65, 84, 259, 264, 256

The dictionary initially contains the ASCII symbols, their indexes are their ASCII codewords.

read	written	added
65	A	
67	C	AC , 257
65	A	CA , 258
71	G	AG , 259
65	A	GA , 260
65	A	AA , 261
84	T	AT , 262
259	AG	TA , 263
264	AGA	AGA , 264
256		

The critical situation occurs when reading the index 264 because, at that moment, no word of the dictionary has this index. □

1.5.3 Implementation

We describe how the dictionary, which is the main data structure of the method, can be implemented. It is natural to consider two implementations adapted for the two phases respectively because the dictionary is not manipulated in the same manner during these phases. They have in common a dictionary implemented as a trie stored in a table D . A node p of the trie is just an index on the table D . It has the three following components:

- $parent(p)$, a link to the parent node of p ,
- $label(p)$, a character,
- $code(p)$, the codeword (index in the dictionary) associated with p .

```

LZW-CODING (fin, fout)
1  count ← -1
2  for each character a ∈ Σ
3      do count ← count + 1
4          HASH-INSERT (D, (-1, a, count))
5  count ← count + 1
6  HASH-INSERT (D, (-1, END, count))
7  p ← -1
8  while not end of file fin and a is the next symbol
9      q ← HASH-SEARCH (D, (p, a))
10     if q = NIL
11         then write code(p) on 1 + log(count) bits in fout
12             count ← count + 1
13             HASH-INSERT (D, (p, a, count))
14             p ← HASH-SEARCH (D, (-1, a))
15         else p ← q
16 write p on 1 + log(count) bits in fout
17 write code(HASH-SEARCH (D, (-1, END))) in 1 + log(count) bits in fout

```

Figure 1.24: LZW encoding algorithm.

In the compression algorithm shown in Figure 1.24, for a node p we need to compute its child according to some letter a . This is done with by hashing, with a hashing function defined on pairs in the form (p, a) . This provides a fast access to the children of a node.

The function HASH-SEARCH, with input $(D, (p, a))$, returns the node q such that $parent(q) = p$ and $label(q) = a$, if such a node exists and NIL otherwise. The procedure HASH-INSERT, with input $(D, (p, a, c))$, inserts a new node q in the dictionary D with $parent(q) = p$, $label(q) = a$ and $code(q) = c$.

For the decompression algorithm, no hashing technique is necessary on the table representation of the trie that implements the dictionary. Having the index of the next segment, a bottom-up walk in the trie produces the mirror image of the expected segment. A stack is then used to reverse it. We assume that the function $string(c)$ performs this specific work for a code c . The bottom-up walk follows the parent links of the data structure. The function $first(w)$ gives the first character of the word w . These features are part of the decompression

```

LZW-DECODING (fin, fout)
1  count ← -1
2  for each character a ∈ Σ
3      do count ← count + 1
4          HASH-INSERT (D, (-1, a, count))
5  count ← count + 1
6  HASH-INSERT (D, (-1, END, count))
7  c ← first code on 1 + log(count) bits from fin
8  write string(c) in fout
9  a ← first(string(c))
10 repeat
11     d ← next code on 1 + log(count) from fin
12     if d > count
13         then count ← count + 1
14             parent(count) ← c
15             label(count) ← a
16             write string(c)a in fout
17             c ← d
18     else a ← first(string(d))
19         if a ≠ END
20             then count ← count + 1
21                 parent(count) ← c
22                 label(count) ← a
23                 write string(d) in fout
24                 c ← d
25     else exit
26 forever

```

Figure 1.25: LZW decoding algorithm.

algorithm displayed in Figure 1.25.

The Ziv-Lempel compression and decompression algorithms run both in time linear in the sizes of the files provided the hashing technique is implemented efficiently. Indeed, it is very fast in practice, except when the table becomes full and should be reset. In this situation, usual implementations also reset the whole dictionary to its initial value.

The main advantage of Ziv-Lempel compression method, compared to Huffman coding, is that it captures long repeated segments in the source file and thus often yields better compression ratios.

1.6 Mixing several methods

We describe simple compression methods and then an example of a combination of several of them, basis of the popular **bzip** software.

1.6.1 Run Length Encoding

The aim of Run Length Encoding (RLE) is to efficiently encode repetitions occurring in the input data. Let us assume that it contains a good quantity of repetitions of the form $aa\dots a$ for some character a ($a \in \Sigma$). A repetition of k consecutive occurrences of letter a is replaced by $\&ak$, where the symbol $\&$ is a new character ($\& \notin \Sigma$).

The string $\&ak$ that encodes a repetition of k consecutive occurrences of a is itself encoded on the binary alphabet $\{0, 1\}$. In practice, letters are often represented by their ASCII code. Therefore, the codeword of a letter belongs to $\{0, 1\}^k$ with $k = 7$ or 8 . Generally there is no problem in choosing or encoding the special character $\&$. The integer k of the string $\&ak$ is also encoded on the binary alphabet, but it is not sufficient to translate it by its binary representation, because we would be unable to recover it at decoding time inside the stream of bits. A simple way to cope with this is to encode k by the string $0^\ell \text{bin}(k)$, where $\text{bin}(k)$ is the binary representation of k , and ℓ is the length of it. This works well because the binary representation of k starts with a 1 so there is no ambiguity to recover ℓ by counting during

the decoding phase. The size of the encoding of k is thus roughly $2 \log k$. More sophisticated integer representations are possible, but none is really suitable for the present situation. Simpler solution consists in encoding k on the same number of bits as other symbols, but this bounds values of ℓ and decreases the power of the method.

1.6.2 Move To Front

The Move To Front (MTF) method may be regarded as an extension of Run Length Encoding or a simplification of Ziv–Lempel compression. It is efficient when the occurrences of letters in the input text are localized into relatively short segment of it. The technique is able to capture the proximity between occurrences of symbols and to turn it into a short encoded text.

Letters of the alphabet Σ of the input text are initially stored in a list that is managed dynamically. Letters are represented by their rank in the list, starting from 1, rank that is itself encoded as described above for RLE.

Letters of the input text are processed in an on-line manner. The clue of the method is that each letter is moved to the beginning of the list just after it is translated by the encoding of its rank.

The effect of MTF is to reduce the size of the encoding of a letter that reappears soon after its preceding occurrence.

1.6.3 Integrated example

Most compression software combine several methods to be able to compress efficiently a large range of input data. We present an example of this strategy, implemented by the UNIX command **bzip**.

Let $y = y[0]y[1] \cdots y[n-1]$ be the input text. The k -th rotation (or conjugate) of y , $0 \leq k \leq n-1$, is the string $y_k = y[k]y[k+1] \cdots y[n-1]y[0]y[1] \cdots y[k-1]$.

We define the *BW* transformation as $BW(y) = y[p_0]y[p_1] \cdots y[p_{n-1}]$, where $p_i + 1$ is such that y_{p_i+1} has rank i in the sorted list of all rotations of y .

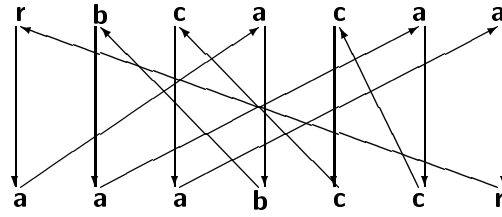


Figure 1.26: Example of text $y = \mathbf{baccara}$. Top line is $BW(y)$ and bottom line the sorted list of letters of it. Top-down arrows correspond to succession of occurrences in y . Each bottom-up arrow links the same occurrence of a letter in y . Arrows starting from equal letters do not cross. The circular path is associated with rotations of the string y . If the starting point is known, the only occurrence of letter \mathbf{b} here, following the path produces the initial string y .

Source texts	French	C sources	Alphabet	Random
Sizes in bytes	62816	684497	530000	70000
Huffman	53.27%	62.10%	72.65%	55.58%
Ziv-Lempel	41.46%	34.16%	2.13%	63.60%
Factor	47.43%	31.86%	0.09%	73.74%

Figure 1.27: Sizes of texts compressed with three algorithms.

It is remarkable that y can be recovered from both $BW(y)$ and a position on it, starting position of the inverse transformation (see Figure 1.26). This is possible due to the following property of the transformation. Assume that $i < j$ and $y[p_i] = y[p_j] = a$. Since $i < j$, the definition implies $y_{p_i+1} < y_{p_j+1}$. Since $y[p_i] = y[p_j]$, transferring the last letters of y_{p_i+1} and y_{p_j+1} to the beginning of these words does not change the inequality. This proves that the two occurrences of a in $BW(y)$ are in the same relative order as in the sorted list of letters of y . Figure 1.26 illustrates the inverse transformation.

Transformation BW obviously does not compress the input text y . But $BW(y)$ is compressed more efficiently with simple methods. This is the strategy applied for the command **bzip**. It is a combination of the BW transformation followed by MTF encoding and RLE encoding. Arithmetic coding, a method providing compression ratios slightly better than Huffman coding, may also be used.

1.7 Experimental results

The table of Figure 1.27 contains a sample of experimental results showing the behavior of compression algorithms on different types of texts. The table is extracted from (Zipstein, 1992).

The source files are: French text, C sources, Alphabet, and Random. Alphabet is a file containing a repetition of the line `abc...zABC...Z`. Random is a file where the symbols have been generated randomly, all with the same probability and independently of each others.

The compression algorithms reported in the table are: the Huffman algorithm of Section 1.2, the Ziv-Lempel algorithm of Section 1.5, and a third algorithm called Factor. This latter algorithm encodes segments of the source text as Ziv-Lempel algorithm does. But the segments are taken among all segments already encountered in the text before the current position. The method gives usually better compression ratio but is more difficult to implement. Compression based on arithmetic as presented in this chapter gives compression ratios slightly better than Huffman coding.

The table of Figure 1.27 gives in percentage the sizes of compressed files. Results obtained by Ziv-Lempel and Factor algorithms are similar. Huffman coding gives the best result for the Random file. Finally, experience shows that exact compression methods often reduce the size of data to 30%–50% of their original size.

Table 1.1 contains a sample of experimental results showing the behavior of compression algorithms on different types of texts from the Calgary Corpus: `bib` (bibliography), `book1` (fiction book), `news` (USENET batch file), `pic` (black and white fax picture), `progC` (source code in C) and `trans` (transcript of terminal session).

The compression algorithms reported in the table are: the Huffman coding algorithm implemented by `pack`, the Ziv-Lempel algorithm implemented by `gzip-b` and the compression based on the *BW* transform implemented by `bzip2-1`.

Additional compression results can be found at <http://corpus.canterbury.ac.nz>.

Sizes in bytes	111,261	768,771	377,109	513,216	39,611	93,695	
Source Texts	bib	book1	news	pic	progc	trans	Average
pack	5.24	4.56	5.23	1.66	5.26	5.58	4.99
gzip-b	2.51	3.25	3.06	0.82	2.68	1.61	2.69
bzip2-1	2.10	2.81	2.85	0.78	2.53	1.53	2.46

Table 1.1: Compression results with three algorithms: Huffman coding (**pack**), Ziv-Lempel coding (**gzip-b**) and Burrows-Wheeler coding (**bzip2-1**). Figures give the number of bits used per character (letter). They show that **pack** is the least efficient method and that **bzip2-1** compresses slightly more than **gzip-b**.

1.8 Research Issues and Summary

The statistical compression algorithm is from Huffman (1951). The UNIX command **pack** implements the algorithm.

The dynamic version was discovered independently by Faller (1973) and Gallager (1978). Practical versions were given by Cormack and Horspool (1984) and Knuth (1985). A precise analysis leading to an improvement was presented in (Vitter, 1987). The command **compact** of UNIX implements the dynamic Huffman coding.

It is unclear to whom precisely should be attributed the idea of data compression using arithmetic coding. It is sometimes refer to Elias (1963), and has become popular after the publication of the article of Witten, Neal, and Cleary. An efficient data structure for the tables of frequencies is due to Fenwick (1994). The main interest in arithmetic coding for text compression is that the two different processes “modeling” the statistics of texts and “coding” can be made independent modules.

Several variants of the Ziv-Lempel algorithm exist. The reader can refer to the books of Bell, Cleary, and Witten (1990) or Storer (1988) for a discussion on them.

The *BW* transform is from Burrows and Wheeler (1994).

The books of Held (1991) and Nelson (1992) present practical implementations of various compression algorithms, while the book of Crochemore and Rytter (2002) overflows the strict topic of text compression and described more algorithms and data structures related to

texts manipulations.

1.9 Defining Terms

Codeword: sequence of bits of a code corresponding to a symbol.

Prefix: a word $u \in \Sigma^*$ is a prefix of a word $w \in \Sigma^*$ if $w = uz$ for some $z \in \Sigma^*$.

Prefix code: set of words such that no word of the set is a prefix of another word contained in the set. A prefix code is represented by a coding tree.

Segment: a word $u \in \Sigma^*$ is a segment of a word $w \in \Sigma^*$ if u occurs in w , i.e. $w = vuz$ for two words $v, z \in \Sigma^*$. (u is also referred to as a factor or a subword of w)

Trie: tree in which edges are labeled by letters or words.

1.10 References

- Cormack, G.V., Horspool, R.N.S. 1984. Algorithms for adaptive Huffman Codes. *Inf. Process. Lett.* 18(3):159–165.
- Faller, N. 1973. An adaptive system for data compression. In *Record of the 7th Asilomar Conference on Circuits, Systems, and Computers*. 593-597.
- Fenwick, P.M. 1994. A new data structure for cumulative frequency tables. *Software—Practice and Experience*. 24(7):327–336.
- Gallager, R.G. 1978. Variations on a theme by Huffman. *IEEE Trans. Inf. Theory*. 24(6):668–674.
- Huffman, D.A. 1951. A method for the construction of minimum redundancy codes. *Proceedings of the I.R.E.* 40:1098–1101.
- Knuth, D.E. 1985. Dynamic Huffman coding. *J. Algorithms*. 6(2):163–180.
- Vitter, J.S. 1987. Design and analysis of dynamic Huffman codes. *J. ACM*. 34(4):825–845.

Welch, T.A. 1984. A technique for high-performance data compression. *IEEE Computer*. 17(6):8-19.

Witten, I.H., Neal, R., Cleary, J.G. 1987. Arithmetic coding for data compression. *Comm. ACM*. 30(6):520-540.

Ziv, J., Lempel, A. 1977. A universal algorithm for sequential data compression. *IEEE Trans. Inf. Theory*. 23(3):337-343.

1.11 Further Information

The set of algorithms presented in this chapter provides the basic methods for data compression. In commercial software they are often combined with other more elementary techniques that are described in textbooks. A wider panorama of data compression algorithms on texts may be found in several books such as:

- Bell, T.C., Cleary, J.G., and Witten, I.H. 1990. *Text Compression*, Prentice Hall, Englewood Cliffs, New Jersey.
- Crochemore, M., and Rytter, W. 2002. *Jewels of Stringology*, World Scientific Publishing.
- Held, G. 1991. *Data Compression*, John Wiley & Sons.
- Nelson, M. 1992. *The Data Compression Book*, M&T Books.
- Storer, J.A. 1988. *Data Compression: Methods and Theory*, Computer Science Press.

Research papers in text data compression are disseminated in a few journals, among which are: *Communications of the ACM*, *Journal of the ACM*, *Theoretical Computer Science*, *Algorithmica*, *Journal of Algorithms*, *Journal of Discrete Algorithms*, *SIAM Journal on Computing*, *IEEE Trans. Information Theory*.

An annual conference presents the latest advances of this field of research:

- *Data Compression Conference*, which is regularly held at Snowbird (Utah) in spring.

Two other conferences on pattern matching also present research issues in this domain:

- *Combinatorial Pattern Matching* (CPM), which started in 1990 and was held in Paris (France), London (England), Tucson (Arizona), Padova (Italy), Asilomar (California), Helsinki (Finland), Laguna Beach (California), Aarhus (Denmark), Piscataway (New Jersey), Warwick University (England), Montreal (Canada), Jerusalem (Israel), Fukuoka (Japan), Morelia, Michocán (Mexico), Istanbul (Turkey), Jeju Island (Korea), Barcelona (Spain), and London, Ontario (Canada).
- *Workshop on String Processing* (WSP), which started in 1993 and was held in Belo Horizonte (Brasil), Valparaiso (Chile), Recife (Brasil), and Valparaiso (Chile) and became *String Processing and Information Retrieval* (SPIRE) in 1998 and was held in Santa Cruz (Bolivia), Cancun (Mexico), A Coruña (Spain), Laguna de San Rafael (Chile), Lisbon (Portugal), Manaus (Brazil), Padova (Italy), Buenos Aires (Argentina), Glasgow (Scotland), and Santiago (Chile),

And general conferences in computer science often have sessions devoted to data compression algorithms.