

A formula based approach to Arithmetic Coding

ARUNDALE RAMANATHAN

Siara Logics (cc)
arun@siara.cc

Abstract

The Arithmetic Coding process involves re-calculation of intervals for each symbol that need to be encoded. This article discovers a formula based approach for calculating compressed codes and provides proof for deriving the formula from the usual approach. A spreadsheet is also provided for verification of the approach. Consequently, the similarities between Arithmetic Coding and Huffman coding are also visually illustrated.

This article presents a formula based approach to Arithmetic Coding. It also explains the mathematical foundation of Arithmetic Coding from a radically different perspective.

This article discovers a formula based approach for calculating compressed codes and provides proof for deriving the formula from the usual approach.

Using a spreadsheet, the new approach is demonstrated by compressing and decompressing a simple string ("Hello World"). Compression using conventional approach is also demonstrated in the same spreadsheet. It can be seen that the same compressed value is obtained using both the methods.

Simply put, the spreadsheet compresses the string "Hello World" having length 11 letters to a 4 byte value 2166290293 (31.298 bits to be more exact).

Radix (logarithmic) is taken as 2 in this article, but any radix can be used.

I Known facts

The basic principles of Arithmetic Coding are explained well in [2].

Referring to section *Example* in this article [3], for a given input string A , with N symbols (letters) and n unique symbols,

- if each unique symbol is represented as a_i , i being the index of symbol after sorting by descending order of weights,

- if each symbol appears k_i times,
- if weights (or probability) of each symbol is given by $w_i = k_i / N$,

we know:

- the length (in bits) of optimal possible code for symbol a_i is $-\log_2(w_i)$ bits, referred as $h(a_i)$.
- the total compressed length L will be $\sum_{i=1}^n -k_i \log_2(w_i)$ or $\sum_{i=1}^n k_i h(a_i)$ bits.

II This work

The formula for $h(a_i)$ has been known for several decades [1]. If this indicates the length of the compressed code (in bits), what is the *value* contained in that length?

Surely it cannot be all 0s or all 1s, in which case the compressed value will simply be 0 or 1. Also, it cannot also take any arbitrary value, as there could be more than one symbol having the same compressed length. So the value is distinct and specific for each symbol. Let us call this *value* as $v(a_i)$.

If the formula for this value could be discovered, it would be simply a matter of concatenating *lengths* of *values* to obtain the code, without having to recalculate the intervals for each symbol, as required in the common approach.

III Formulae

$$v(a_i) = \left(\sum_{j=1}^{i-1} k_j \right) / k_i, \quad \forall i > 1, \quad v(a_1) = 0.$$

$$\text{compressed_value} = \sum_{i=1}^N \left(v(A[i]) / 2^{\sum_{j=1}^i h(A[j])} \right)$$

IV Proof (derivation)

We derive the formula from the common approach, which slices the interval according to the weights w_j . So for coding any symbol a_i , the following value is used:

$$\text{code_value}_i = \sum_{j=1}^{i-1} (w_j * \text{interval_length})$$

$\forall i > 1$. code_value_1 would be 0.

When we use interval as 0 to 1, interval_length is equal to 1, so it becomes:

$$\text{code_value}_i = \sum_{j=1}^{i-1} w_j$$

$\forall i > 1$. code_value_1 would be 0.

However, since we have taken the interval as 0 to 1, the value is a fraction and we have to scale it to get the value we are seeking. The length of this value is $h(a_i)$ bits, so we shift it left as follows to get the desired value:

$$\begin{aligned} v(a_i) &= \left(\sum_{j=1}^{i-1} w_j \right) * 2^{h(a_i)} \\ \Rightarrow v(a_i) &= \left(\sum_{j=1}^{i-1} w_j \right) * 2^{-\log_2(w_i)} \\ \Rightarrow v(a_i) &= \left(\sum_{j=1}^{i-1} w_j \right) * 2^{\log_2(1/w_i)} \\ \Rightarrow v(a_i) &= \left(\sum_{j=1}^{i-1} w_j \right) * (1/w_i) \\ \Rightarrow v(a_i) &= \left(\sum_{j=1}^{i-1} w_j \right) / w_i \\ \Rightarrow v(a_i) &= \left(\sum_{j=1}^{i-1} (k_j/N) \right) / (k_i/N) \end{aligned}$$

$$\therefore v(a_i) = \left(\sum_{j=1}^{i-1} k_j \right) / k_i$$

For all above statements, $i > 1$ and $v(a_1)$ is always 0.

For the above derivation, any interval could be used as long as $h(a_i)$ is then *right aligned* to 0. For example, if we use an interval of 0 to 4294967296 (2^{32}), the value would have to be shifted right by $32 - h(a_i)$.

If interval_end is expressed in terms of powers of 2, say 2^x , then the value should be *right shifted* by $x - h(a_i)$. Note that when interval is 0 to 1, $\text{interval_end} = 2^0$ ($x=0$).

V Application to other coding methods

The same formulae and approach are applicable to other coding methods such as Huffman coding, Shannon-Fano coding where length and value are available. The only difference is that length would be a natural number in these cases.

Once the codes for symbols are obtained using the respective methods, the Frequencies need to be re(verse)-calculated according to the code lengths ($k_i = 2^{-h(a_i)} * N$). Then, the formulas can be applied to obtain the compressed value. This is shown in a separate sheet (Huffman_coding).

A picture for visual comparison between Arithmetic Coding and Huffman Coding is given under the example section below (Fig. 3).

VI Example

Given $A = \text{"Hello World"}$, then

- $N = 11, n = 8,$
- $a_1 = 'l', a_2 = 'o', a_3 = 'H', a_4 = 'e', a_5 = ' ', a_6 = 'W', a_7 = 'r', a_8 = 'd',$ and
- $k_1 = 3, k_2 = 2, k_3 \text{ to } k_8 = 1$
- $w_1 = 0.2727 (3/11), w_2 = 0.1818 (2/11), w_3 \text{ to } w_8 = 0.0909 (1/11)$

then

- $h(a_1) = 1.8745, h(a_2) = 2.4594, h(a_3)$ to $h(a_8) = 3.4594$, and
- $L = 31.2989$

which means after compression, 11 bytes will become 31.2989 bits (around 4 bytes). By applying the formulas, we get:

- $v(a_1) = 0, v(a_2) = 1.5, v(a_3) = 5, v(a_4) = 6, v(a_5) = 7, v(a_6) = 8, v(a_7) = 9, v(a_8) = 10$

- $compressed_value = 2166290392.64712$
(0.50437878646122 unscaled)

The detailed calculations can be seen in the spreadsheet (Fig. 2).

The following picture visually shows placement of letters in compressed value for both Arithmetic and Huffman coding:

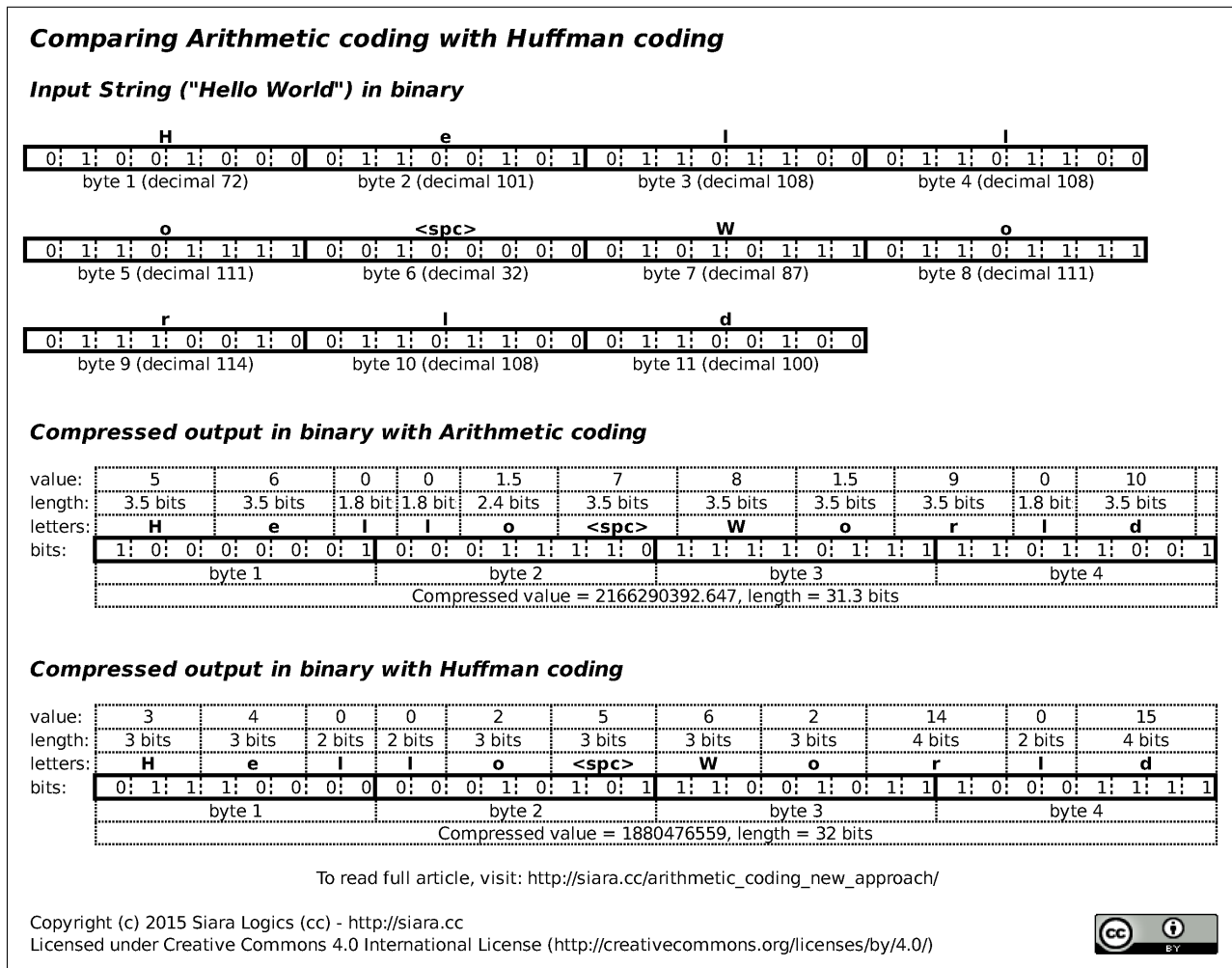


Figure 1: Visual indication of symbol positions in compressed value

The values are the same as those shown in the example spreadsheets. A screenshot of the spreadsheets are also given in the Appendix.

VII Conclusion

The current work simplifies the encoding process and explains Arithmetic coding in simpler terms. However, for practical implementations, the following points need to be considered:

- The formula based approach would heavily depend on the performance of *exp2* function. It is likely to be slower than calculating the interval table for each symbol. Very little research has been done on this aspect.
- While the interval based approach does not allow parallel processing [1], the formula based approach would allow parallel processing.

Till the answers to the above points are available, this work presently serves the following purposes:

- Understand Arithmetic Coding from a different perspective

- Visualize positions of compressed symbols
- Visually compare Arithmetic coding and Huffman coding
- Precursor for further research on entropy coding

References

- [1] Paul G. Howard and Jeffrey Scott Vitter, *Practical Implementations of Arithmetic Coding*, Brown University, April 1992.
- [2] Wikipedia, *Arithmetic Coding*, https://en.wikipedia.org/wiki/Arithmetic_coding, September 2015.
- [3] Wikipedia, *Huffman Coding*, https://en.wikipedia.org/wiki/Huffman_coding, August 2015.

VIII Appendix

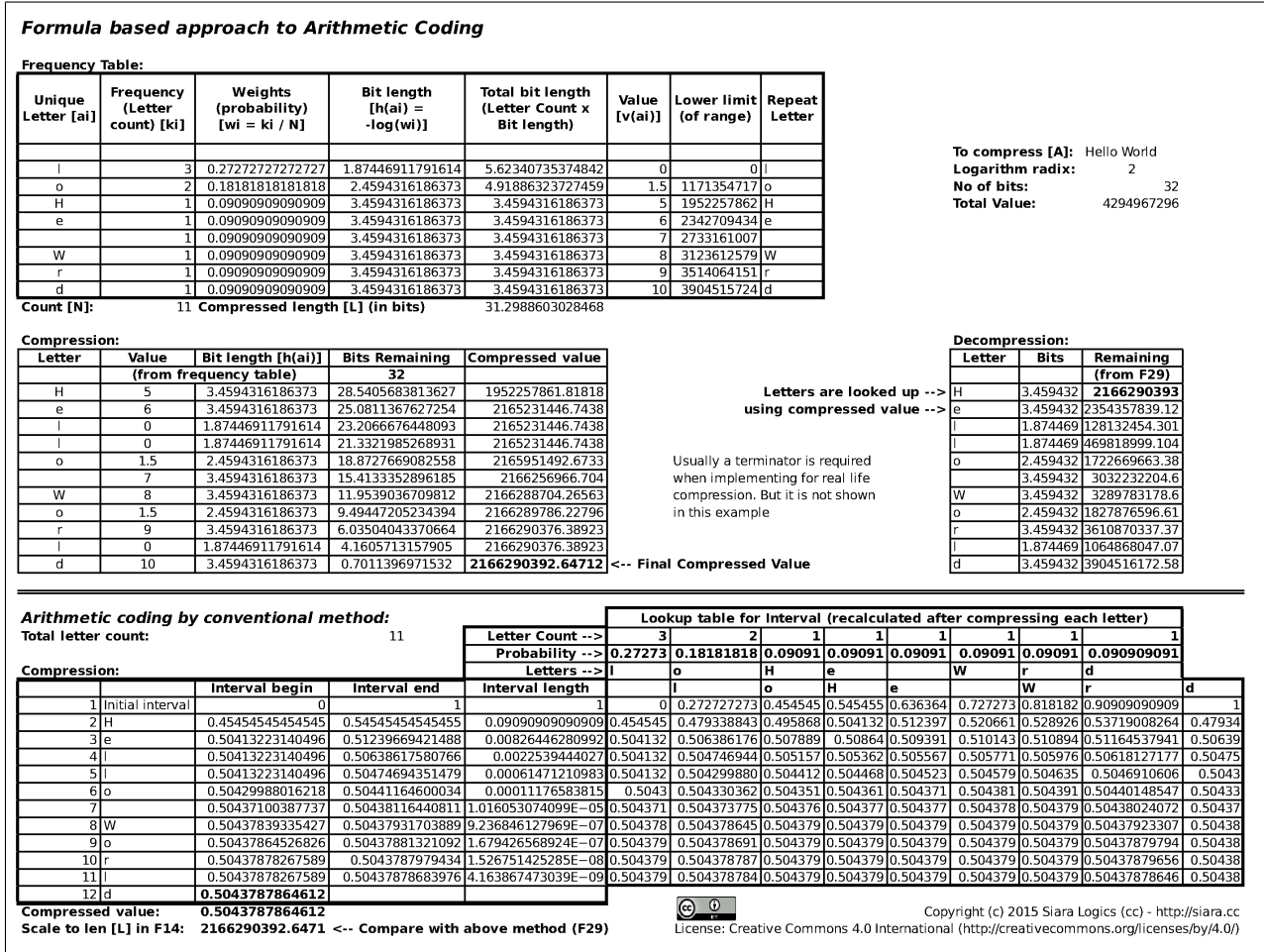


Figure 2: Screenshot of spreadsheet that demonstrates both methods of Arithmetic Coding

Comparison with Huffman Coding

Frequency Table:

Unique Letter [ai]	Frequency (Letter count) [ki]	Weights (probability) [wi = ki / N]	Bit length [h(ai)] (copied from binary tree built elsewhere)	Total bit length (Letter Count x Bit length)	Value [v(ai)]	Lower limit (of range)	Repeat Letter	Frequency [ki] re-calculated from h(ai)	Huffman Code
l	3	0.27272727273	2	6	0	0	l	2.75	00
o	2	0.18181818182	3	6	2	1073741824	o	1.375	010
H	1	0.09090909091	3	3	3	1610612736	H	1.375	011
e	1	0.09090909091	3	3	4	2147483648	e	1.375	100
	1	0.09090909091	3	3	5	2684354560		1.375	101
W	1	0.09090909091	3	3	6	3221225472	W	1.375	110
r	1	0.09090909091	4	4	14	3758096384	r	0.6875	1110
d	1	0.09090909091	4	4	15	4026531840	d	0.6875	1111

Count [N]: 11 Compressed len [L] (in bits): 32 11

Compression:

Letter	Value	Bit length [h(ai)]	Bits Remaining	Compressed value
		(from frequency table)	32	
H	3	3	29	1610612736
e	4	3	26	1879048192
l	0	2	24	1879048192
l	0	2	22	1879048192
o	2	3	19	1880096768
	5	3	16	1880424448
W	6	3	13	1880473600
o	2	3	10	1880475648
r	14	4	6	1880476544
l	0	2	4	1880476544
d	15	4	0	1880476559 <-- Final Compressed Value

Letters are looked up -->
using compressed value -->

Usually a terminator is required when implementing for real life compression. But it is not shown in this example

Decompression:

Letter	Bits	Remaining (from F29)
H	3	1880476559
e	3	2158910584
l	2	91415488
l	2	365661952
o	3	1462647808
	3	3111247872
W	3	3415146496
o	3	1551368192
r	4	3821010944
l	2	1006632960
d	4	4026531840

Given string [A]: Hello World
 Logarithm Radix: 2
 No of bits: 32
 Total Value: 4294967296

Copyright (c) 2015 Siara Logics (cc) - <http://siara.cc>
 Licensed under Creative Commons 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>)

Figure 3: Screenshot of spreadsheet that demonstrates compression using Huffman codes by the same formula