

Optimised Cost Considering Huffman Code for Biological Data Compression

Abstract

Classical Huffman code has widely been used to compress biological datasets. Though a considerable reduction of size of data is obtained by classical Huffman code, yet, more efficient encoding is possible by treating binary bits differently considering requirement of transmission time, energy consumption, and likewise. A number of techniques already modified Huffman code algorithm to obtain optimal prefix-codes for unequal letter costs in order to reduce overall transmission cost (time). In this paper we propose a new approach to improve compression performance of one such extension, cost considering Huffman code, by applying genetic algorithm for optimal allocation of the codewords to the symbols. The approach start with generating a set of codewords for different input symbols based on their frequencies by considering cost (length) of 0 and 1 as α and β , $\alpha < \beta$. This step ensures optimal cost encoding, however, compression performance is not higher than classical Huffman code. To improve compression performance, the second part of the approach uses genetic algorithm to find an optimal allocation of generated codes to different symbols. The idea is to sacrifice some cost to minimise total number of bits, hence, the genetic algorithm works by giving penalty on cost. The performance of the approach is evaluated by applying it to compress some standard biological dataset and comparing the result with the results produced by classical Huffman code and standard cost considering Huffman code.

Keywords: Data Transmission , Huffman code, Cost Considering, Genetic algorithm, Biological data.

1. Introduction

In the recent years, application of battery-powered portable devices, e.g. laptop computers and mobile phones has increased significantly. Proper representation of digital data and their transmission efficiency has become a

primary concern for digital community because it affects the performance, reliability, and the cost of computation in both portable and non-portable devices. CMOS technologies were developed in order to reduce the power consumption both in data processing and transmission. In order to increase transmission speed and reduce transmission cost, parallel data transmission methods are widely used. However, parallel transmission is limited to short distance communications, e.g. locally connected devices, internal buses. Ruling out the possible availability of parallel transmission links over long distance, we are left with its serial alternative only. If we attempt to transfer big files, e.g. DNA sequences, over a serial transmission link then it would take a significant amount of time. However, we cannot overlook this problem because at present parallel processing is widely used to increase throughput and in parallel processing architecture, processing units are usually distributed in different physical locations and task sharing is a must in such architecture.

Data encoding techniques came into action to improve the data transmission efficiency over serial communication medium by compressing data before transmitting. Efficiency can be measured in terms of incurred cost, required storage space, consumed power, time spent and likewise. Data must be encoded to meet the purposes like: unambiguous retrieval of information, efficient storage, efficient transmission and etc. Let a message consist of sequences of characters taken from an alphabet Σ , where $\alpha_1, \alpha_2, \alpha_3 \dots, \alpha_r$ are the elements that represent the characters in the source Σ . The length of α_i represents its cost or transmission time, i.e., $c(\alpha_i) = \text{length}(\alpha_i)$. A codeword w_i is a string of characters in Σ , i.e., $w_i \in \Sigma^+$. If a codeword is $w_i = \alpha_{i1}, \alpha_{i2}, \dots, \alpha_{in}$, then the length or cost of the codeword is the sum of the lengths of its constituent elements:

$$\text{cost}(w_i) = \sum_{j=1}^n c(\alpha_{ij}) \quad (1)$$

If all the elements of a codeword has unit cost or length then the cost of the codeword is equivalent to the length of the codeword. However, it is not necessary for the elements in the codeword to have equal length or cost. For example, in Morse Code all the ASCII characters are encoded as sequence of dots (·) and dashes (–) where a dash is three times longer than a dot in duration [1]. However, the Morse code scheme suffers from the prefix problem [2]. Ignoring the prefix problem, Morse code results in a tremendous savings of bits over ASCII representation. Using Morse code, we can treat

the binary bits differently; 0 as a dot and 1 as a dash. Even if we consider the voltage level to represent the binary digits then they are still different. Table 1 shows the logic level to represent binary digits in CMOS and TTL technologies.

Table 1: Example of binary logic level

Technology	0	1	Notes
CMOS	0 V to $\frac{V_{DD}}{2}$	$\frac{V_{DD}}{2}$ to V_{DD}	V_{DD} = supply voltage
TTL	0 V to 0.8 V	2 V to V_{CC}	V_{CC} is 4.75 V to 5.25 V

As the unequal letter cost problem is not new therefore it has been addressed by different researchers. The more general case where the costs of the letters as well as the probabilities of the words are arbitrarily specified was treated in [3]. A number of other researchers have focused on uniform sources and developed algorithm for the unequal letter costs encoding [4, 5, 6, 7, 8]. Let p_1, p_2, \dots, p_n be the probabilities with which the source symbols occur in a message and the codewords representing the source symbols are w_1, w_2, \dots, w_n then the cost of the code W is:

$$C(W) = \sum_{i=1}^n \text{cost}(w_i) \cdot p_i \quad (2)$$

The aim of producing an optimal code with unequal letter cost is to find a codeword W that consists of n prefix code letters each with minimum cost c_i that produces the overall minimum cost $C(W)$, given that costs $0 < c_1 \leq c_2 \leq c_3 \leq \dots \leq c_n$, and probabilities $p_1 \geq p_2 \geq \dots \geq p_n > 0$.

Huffman code is an efficient data compression scheme that takes into account the probabilities at which different quantization levels are likely to occur and results in fewer data bits on the average. It is widely used to compress biological data, however, all the techniques use the classical form of the Huffman code where bits are treated equally. Out of many variations of the Huffman code where cost of bits are treated unequally, the most recent approach is described in [9]. This approach treated binary bit 0 as a dot (·) and 1 as a dash (–) like Morse code and reduces the transmission cost (time) significantly. Like other variations of the cost considering Huffman code, the compression performance (in terms of number of bits require to encode a

message) of this approach is not also better than classical Huffman code. In this paper, we have proposed a new optimised cost considering Huffman code based on the approach shown in [9]. This new approach optimised the number of bits require to encode a message while treating the binary bits unequally. The efficiency of the method is evaluated by applying it to compress some standard biological dataset.

The rest of the paper is organised as follows: Section 2 presents the background study of the issues in biological data processing and the Huffman code. The proposed approach is described in Section 3. Experimental results and discussion are presented in Section 4 . Finally, concluding remarks are presented in Section 5.

2. Background

2.1. Issues on Biological Data Transmission

The size of biological data including DNA sequences increase with an ever expanding rate and will be bigger and bigger in the future. These Biological data are stored in biology database, the exponential growth of these database become a big problem to all biological data processing methods [10]. Different operation will be applied to these data such as, searching [11], e-mail attachment [12], alignment [13], and transmission on distributed computing [14]. Interestingly, biological data compression can play a key role in all biological data processing.

A recent deluge of interest in the development of new tools for biological data processing, these all algorithms needs an efficient methods for data compression. The main objective of data compression methods is minimizing the number of bits in the data representation. In [15] authors propose a new general data structure and data encoding approach for the efficient storage of genomic data. This method encode only the differences between a genome sequence and a reference sequence, the method use different encoding scheme from fixed codes such as Golomb, Elias codes and variable codes such as Huffman codes. Other methods based on same idea to encode only the difference between reference sequence and the target one, Authors in [12] uses Huffman code for encoding difference between sequence to sent it as an email attachment, but these methods suffer that they must sent the reference sequence for at least one time for each species.

Wang and Zhang [16] proposed a new scheme for referential compression of genomes based on the chromosome level. The Algorithm aim to search for

longest common subsequence between matching parts and the differences encoded using Huffman coding.

All previous studies focus only on the differences and the relation between continuation of the sequence, and without improvement of the encoding scheme.

2.2. Huffman Codes

In computer science and information theory, Huffman code is an entropy encoding algorithm used for lossless data compression. It takes into account the probabilities at which different symbols are likely to occur and results into fewer data bits on the average. For any given set of symbols and associated occurrence probabilities, there is an optimal encoding rule that minimises the number of bits needed to represent the source. Encoding symbols in predefined fixed length code, does not attain an optimum performance, because every character consumes equal number of bits irrespective to their degree of contribution to the whole message. Huffman code tackles this by generating variable length codes, given a probability usage frequency for a set of symbols. It generates prefix-code to facilitate unambiguous retrieval of information. A scheme of prefix code assigns codes to letters in Σ to form codeword w_i such that none of them is a prefix to another. For example, the codes $\{1, 01, 001, 0001\}$ and $\{000, 001, 011, 111\}$ are prefix-free, whereas the code $\{1, 01, 100\}$ is not, because 1 is a prefix in 100.

Applications of Huffman code are pervasive throughout computer science. The algorithm to completely perform Huffman encoding and decoding is explained by [17]. It can be used effectively where there is a need for a compact code to represent a long series of a relatively small number of distinct bytes. For example, Table 1 shows 8 different ASCII characters, their frequencies, ASCII codes and the codewords generated for those symbols using Huffman code. It is seen from the table that the codeword to represent each character is compressed and the most frequent character gets the shortest code. In this example, the compression ratio obtained by Huffman code is 64.16%.

There are many other variants of Huffman codes that compress source data to reduce data size and/or transmission cost. For example, Mannan and Kaykobad introduced block technique in Huffman coding which overcomes the limitation of reading whole message prior to encoding[18]. In classical Huffman coding scheme, the letter costs are considered as equal. The unequal letter cost versions of Huffman codes scheme are proposed in [19, 20, 21, 22]. In the unequal letter cost version of the classical Huffman code, letters of the

Table 2: Example of application of Huffman Code to compress ASCII characters

Symbols	Frequency	ASCII Code	Codewords using Huffman Code
A	50	01000001	00
B	35	01000010	101
C	42	01000011	110
D	22	01000100	1001
E	65	01000101	01
F	25	01000110	1111
G	9	01000111	1000
H	23	01001000	1110

alphabet are considered as unequal. Recently, in [9] a method is proposed to show the effects of unequal bits cost on classical Huffman code. The idea of this method is to assign the most frequent symbol the minimum cost and the least frequent symbol the maximum cost code, whereas classical Huffman code assigns most frequent symbol the minimum length and the least frequent symbol the maximum length code.

3. Approach

3.1. Proposed Scheme

A genome is a stretch of DNA (or RNA) that codes for a polypeptide (protein), that is a set of amino acids bound together in specific order. Each genomic sequence consist of nucleotide bound together, which are interpreted by the cellular machinery in groups of three, called triplets []. This is the main raison to divide the whole sequence in a set of triplet and give a code to each triplet. The first step in the optimised cost considering algorithm is cutting the sequence in triplet, then compute the frequency of each triplet in the whole sequence. This table of frequencies are used by the cost considering Huffman code to generate low cost code to each triplet (frequency). Finally these codes with frequencies are used by the optimised cost considering algorithm to generate the optimal allocation with a given penalty.

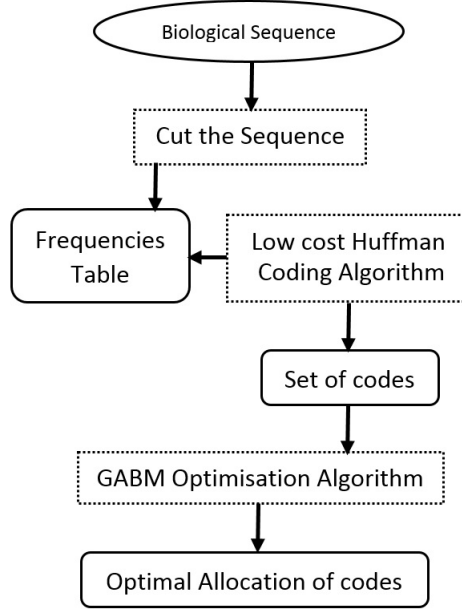


Figure 1: The proposed Scheme

3.2. Cost Considering Huffman code

The classical Huffman algorithm aims at reducing total number of bits and it constructs the tree in a bottom up fashion. It is shown in [23] that if the costs of letters are considered unequal then the straightforward bottom up greedy approach does not work. Authors in [9] uses a top down approach to build a binary tree considering unequal letter cost of bits. They considered cost (length) of 0 and 1 as integer constants α and β , and $\alpha < \beta$. The complete algorithm to obtain an optimal prefix-free code for unequal letter cost is shown below. The input to the algorithm are the distinct symbols contained in the message to be encoded and their frequencies. The process of creating the binary tree starts with a single node (root node) and it is initialised with cost 0. After that, two child (leaf) nodes are created for the root node, i.e. level of the tree is increased by one. Cost of the left child is calculated as the summation of the cost of its parent node and the length of the left arc, and cost of the right child is calculated as the summation of the cost of its parent node and the length of the right arc. Length of left and right arcs are actually the cost (length) of 0 (α) and 1 (β) respectively. The next step is to take a child node with least cost and create two child nodes for

it and make it a parent node. In this way, in every iteration the child node with least cost becomes a parent node with two new child nodes. Creation of new child nodes is stopped when total number of child nodes become equal to the number of distinct symbols needed to be encoded.

Now the tree T is constructed and the cost of the tree actually depends on how we assign the frequencies to the leaf nodes. The overall cost will be minimised if the leaves with highest cost always have smaller or equal weight (frequency). To fulfil this condition the leaves of the T are enumerated in non-decreasing order of their cost, i.e., $cost(l_1) \leq cost(l_2) \leq \dots \leq cost(l_n)$, and that $f_1 \geq f_2 \geq \dots \geq f_n$, where l_i and f_i are leaf node and frequency of distinct symbol respectively for $i = 1, 2, \dots, n$. The frequency or weight of parent nodes are calculated as the sum of its child nodes, and it continues upwards until the root node is reached. After that, the algorithm checks for any possible conflicts between all pair of nodes. Two nodes are considered as conflicted if the node with higher cost has higher frequency violating the above condition, i.e., if $cost(l_i) > cost(l_j)$ and $f(l_i) > f(l_j)$, then there remains a conflict. If there remains a conflict between nodes, then it is resolved by swapping the nodes and recalculating the cost of the tree downward and frequency of the nodes upward. When all the conflicts if existed are resolved then the algorithm generates codes for each of the distinct symbols.

3.3. Optimisation of the Codes

3.3.1. Problem formulation

The problem of finding the best allocation of codes to each symbol can be modelled as an Assignment Problems with Constraint, the problem is formulated as follows :

Definition: Given a set of codes $C = \{C_1, C_2 \dots C_n\}$, and a set of frequencies $C = \{Q_1, Q_2 \dots Q_n\}$. For each code we have the length of the code $|C_i|$ (number of bits) and the cost of the code S_{C_i} (cost of ones and zeros), the objective is to assign to each frequency a code in order to minimize the total number of bits, while respecting the initial assignment total cost S_t with a given penalty λ . This penalty coefficient represent the allowed Marge can be sacrificed for cost to optimise number of bits

The Objective Function is :

$$Minimise \sum (|C_i| \times Q_j) \quad (3)$$

Algorithm 1 Cost-considering / Unequal bit cost Coding

Require: Distinct symbols contained in the message to be encoded and their frequencies

Ensure: Non-uniform / variable letter cost i.e, Cost-considering balanced tree

```
1: for each distinct symbol  $i$  do
2:    $Enqueue(max\_Q, frequency [ i ])$ 
3: end for
4: create a root node
5:  $cost [ root ] \leftarrow 0$ 
6:  $Enqueue(min\_Q, cost [ root ])$ 
7: Define costs of the left and right child of the binary tree
8: repeat
9:    $cost\_of\_parent\_node \leftarrow Dequeue(min\_Q)$ 
10:  create  $left$  and  $right$  child for this node
11:   $cost [ left\_child ] \leftarrow cost\_of\_parent\_node + left\_child\_cost$ 
12:   $Enqueue(min\_Q, cost [ left\_child ])$ 
13:   $cost [ right\_child ] \leftarrow cost\_of\_parent\_node + right\_child\_cost$ 
14:   $Enqueue(min\_Q, cost [ right\_child ])$ 
15:  Mark parent node as explored
16: until  $2(n - 1)$  nodes are created
17: while  $min\_Q \neq \emptyset$  do
18:    $leaf\_node \leftarrow Dequeue(min\_Q)$ 
19:    $frequency[leaf\_node] \leftarrow Dequeue(max\_Q)$ 
20: end while
21: for each parent node  $j$  do
22:    $frequency [ j ] \leftarrow frequency [ left\_child ] + frequency [ right\_child ]$ 
23: end for
24: repeat
25:   if conflict between nodes then
26:     resolve conflict by swapping conflicted nodes
27:     calculate and reassign cost of all affected nodes
28:     calculate and reassign frequency of all affected nodes
29:   end if
30: until all conflicts are resolved
```

while :

$$\sum(|S_{C_i}| \times Q_j) \leq (\lambda + 1)S_t \quad (4)$$

3.3.2. Basic Genetic Algorithm

Genetic Algorithm (GA) is a bio-inspired meta-heuristics algorithm developed by [24]. GA is a stochastic optimisation algorithm imitate the natural evolution process of genomes. GA started by generate a population of random feasible solutions. The optimisation process of GA take this generate at each iteration a new population based on the previous one. The process can be described as follow: Firstly the selection of two or more solution among the current population, by one of the well-known selection techniques [25]. These selected solutions will be considered as a parents, the second operation of the genetic algorithm is the crossover, which take these parents as inputs to generate other new solutions considered as sons. The third operator of the genetic algorithm is the mutation which ensure a good diversification on the search process. The new solutions can be mutate according to a given mutation probability. The mutation change a position or more on the solution with new value different to the previous one. The quality of each solution is computed with the fitness function which control the evolution of the GA population by the deletion of the worst solution and insertion of the good solutions among parents and sons. This processes is repeated until the stopped criteria is achieved which can be the number of generation or if the population is stabilized.

3.3.3. GA for Bits minimisation

The main objective of the GA optimisation algorithm for bits minimisation (GaBm) problem is to assign to each frequency a specific code. The encoding of an optimisation problem solution into a chromosome is one of the most important issue to obtaining a good optimisation results. GaBM uses a two array fixed length to "64" which is the number of combination for all nucleotides, the first contain the frequencies of each codes and the second contain the cost of each codes. In this way, our genetic algorithm will work with the two array and uses the entry index on the allocation process. The genetic algorithm are stochastic algorithm based on random evolution. Generally the initial population is generated in a random affectation. In the GaBm algorithm the population contain firstly the affectation given by the cost considering Huffman code algorithm and secondly the rest of the population is random generated, but all generated solution must satisfy the

initial cost constraints (step 1). The evolution of the population is the key of the optimisation algorithm of genetic algorithm. During each generation the process start with the selection of a proportion of the population to breed a new generation. In the literature many selection methods have been proposed to guide the population evolution [25]. The existent selection methods are varied from a random selection to heuristic based selection. we have chosen to select randomly the part of the population to be processed, as the heuristic methods are very time-consuming (step 3). After that the operations of genetic algorithms are applied for the initial population to generate a new generation of the population (see figure 2). Firstly the crossover operation are applied to these two selected solution (considered as parents) to generate two new solutions (considered as sons)(step 4). In the literature many crossover techniques have been used in genetic algorithm [26], such as one-point crossover which divide the chromosome on two fragments and recombine the second fragments by the other chromosome second fragment, two-point crossover which divide the chromosome on three fragment and recombine the middle fragment by the middle fragments on the other chromosome, and many other crossover techniques to allow a good convergence of the algorithm, in our case we have used the two point crossover with two parameters, the first parameter α is a random value in $[0,63]$, which represent the first cut point and β is a random value too in $[0,63]$, which represent the number of position to be crossed. we used a two random parameters to ensure a good diversification on the whole search space (step 4). These two new generated children may contain conflict like finding a duplicated code allocated to two different frequencies in the solution, so a regulation step is done to ensure the correctness of the solution (see figure 2). Secondly these two new solutions are mutated according to a predefined probability γ , the best value of mutation rate is very problem specific (step 5). The probability γ are fixed to 0.2 to explorer a few position on the solution. The mutation operator used to maintain genetic diversity from one generation of a population of genetic algorithm chromosomes to the next. In our case we have used the mutation as a random swap mutation operator (see figure 3). Each new generated solution must satisfy the cost constraint, these children must be valid with satisfying the cost constraint. The next step is to add these two new solutions (children) to the population (step 7) (see figure 3). Finally the new population are ranked by fitness (step 8), and the worst solution are deleted until the initial size of the population are achieved (step 9). the whole process are repeated until the max number of operation is achieved

(step 10).

Algorithm 2 GA for bits minimisation

- 1: Population initialization (P).
 - 2: **repeat**
 - 3: Select two solutions S_1, S_2 form P.
 - 4: Crossover S_1, S_2 to generate S_{11}, S_{21} (Children).
 - 5: Mutate S_{11}, S_{21} .
 - 6: Validate children with cost constraint (equation 4).
 - 7: Add children to population
 - 8: Rank the population by fitness
 - 9: Remove worst candidates until population limit
 - 10: **until** Max number of generation not achieved
 - 11: Display the best solution from the population P;
-

4. Results And Discussion

The approach has been evaluated with different real genomic biological data, these genomes were downloaded from a recent version of The National Center for Biotechnology Information (NCBI) available on ([http : //www.ncbi.nlm.nih.gov](http://www.ncbi.nlm.nih.gov)). We focused on the sequences alone, ignoring any header and any other exogenous information. In table () the different data set are described with the size of each of them and the references on the biological data bank.

The table 3 present the results founded by the basic Huffman code, cost considering algorithm and optimised cost considering algorithm. The result show that the number of bits of Huffman algorithm is the minimum number among the other algorithm but the cost is very high. The cost considering algorithm improve the quality of the generated codes in terms of cost but the number of bits. The optimised cost considering algorithm try to find the best allocation of codes to frequencies while the cost constraint respected. The table 4 present the best founded number of bits with different penalty, for each genome we find the max number of useful penalty, after this value, increasing the penalty are in-useful (See figure 3), the number of bits achieve the minimum number but the cost stop decreasing (point 1 in figure 3) and this number of bits stabilized while the cost still increasing until it stabilized

Table 3: Dataset description

Data sets	Name	Size (MegaOctet)	Reference
Genome 1	Mycobacterium smegmatis	6,988,302	CP009496
Genome 2	Amycolatopsis benzoatilytica	8,704,271	NZ_KB912942
Genome 3	Mycobacterium rhodesiae NBB3	6,415,739	CP003169
Genome 4	Streptomyces bottropensis ATCC 25435	8,955,726	NZ_KB911581
Genome 5	Mycobacterium smegmatis str. MC2 155	6,988,269	CP009494
Genome 6	Mycobacterium smegmatis MKD8	7,092,137	NZ_KI421511
Genome 7	Bradyrhizobium WSM471	7,784,016	NZ_CM001442
Genome 8	Amycolatopsis thermoflava N1165	8,677,910	NZ_CM001442
Genome 9	Bacillus thuringiensis Bt407	6,026,843	NZ_CM000747
Genome 10	Bacillus thuringiensis serovar thuringiensis	6,323,123	NZ_CM000748
Genome 11	Pseudomonas aeruginosa 9BR	6,801,503	NZ_AFXI01000001
Genome 12	Bacillus thuringiensis serovar berliner ATCC	6,260,142	NZ_CM000753
Genome 13	Bacillus thuringiensis serovar pakistani	6,037,513	NZ_CM000750
Genome 14	Pseudomonas aeruginosa LES400	6,591,121	CP006982

Table 4: GABM for comparison for cost and number of bits of different approaches with different datasets

Data sets	Huffman Algorithm		CCA		OCCA	
	Cost	Bits	Cost	Bits	Cost	Bits
Genome 1	76787151	37256819	67416213	41332479	67416213	40503061
Genome 2	100425402	48778740	88430665	54662727	88430665	53341739
Genome 3	75940155	36860555	66745619	40927187	66745619	40182241
Genome 4	103552729	50047265	90821835	56147481	90821835	54319891
Genome 5	82234926	39772838	71963876	44254676	71963876	43392820
Genome 6	83454842	40370894	73038795	44864009	73038795	44003571
Genome 7	92539488	44977876	81416359	49790957	81416359	49127335
Genome 8	99613856	48183688	87102639	53817307	87102639	52858893
Genome 9	71876739	34874617	62998800	38469200	62998800 8	37831964
Genome 10	75324432	36576034	66084958	40409906	66084958	39739522
Genome 11	80766360	39092450	70620666	43119140	70620666	42359156
Genome 12	74560825	36179579	65359604	39926210	65359604	39221148
Genome 13	71562941	34737083	62758225	38299629	62758225	37809443
Genome 14	78261299	37843793	68354090	41816580	68354090	41050420

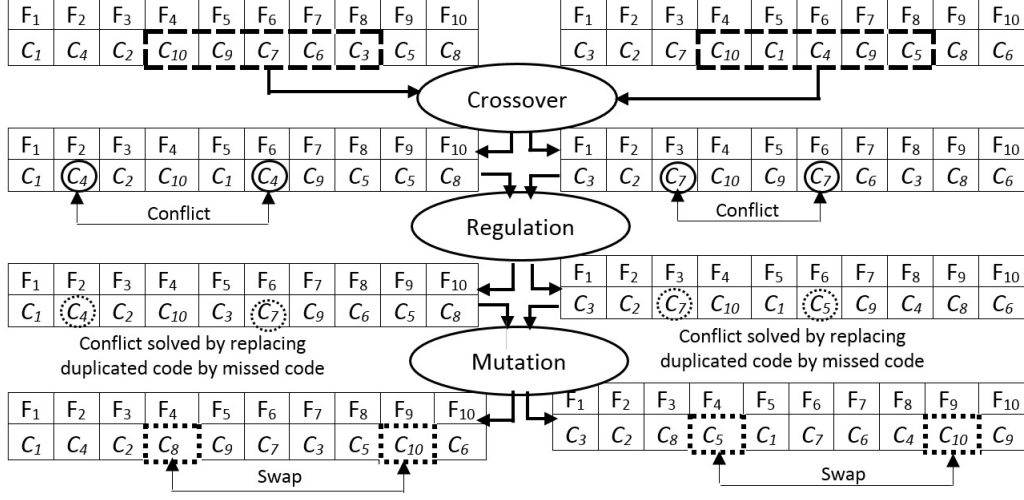


Figure 2: Operations of genetic algorithm

also (point 2 figure 3), after this point the cost and number of bits are stabilized.

5. Conclusion

In this paper we have proposed a new approach for efficient data compression using Huffman code and optimised strategy. The new approach is divided into phases, firstly a cost considering Huffman algorithm are proposed which reduce the cost of the generated codes, these codes are secondary passed by the optimisation algorithm to reduce the global number of bits using a cost penalty.

The proposed approach is tested with biological genomic sequence and a performance comparison is made with the standard Huffman code and the cost considering without optimisation. Simulation results showed that the proposed approach is more robust and efficient compared to other competing algorithms because its penalty based optimisation strategy to search the best allocation of codes to different frequencies.

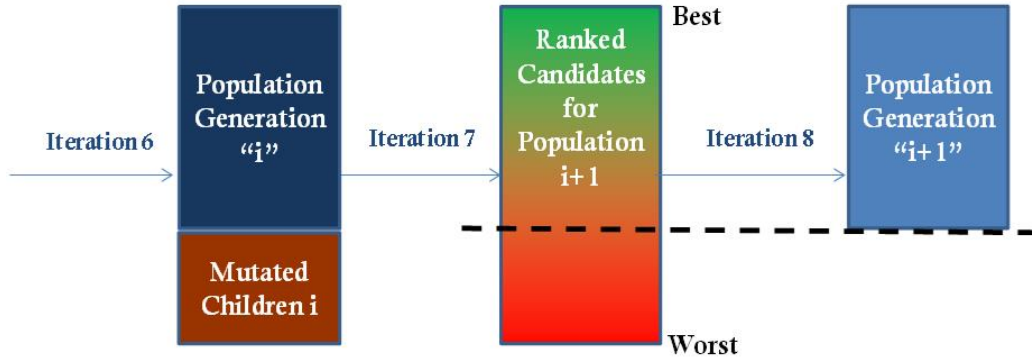


Figure 3: Population Update for genetic algorithm

References

- [1] W. A. Redmond, International morse code, Microsoft Encarta 2009 [DVD] (1964) 275–278.
- [2] P. D. Grunwald, P. M. B. Vitany, Kolmogorov complexity and information theory, *Journal of Logic, Language and Information* 12 (2003) 497–529.
- [3] R. Karp, Minimum-redundancy coding for the discrete noiseless channel, *IRE Transactions on Information Theory* 7 (1) (1961) 27–38.
- [4] E. N. Gilbert, Coding with digits of unequal costs, *IEEE Transactions on Information Theory* 41.
- [5] R. M. Krause, Channels which transmit letters of unequal duration, *Information Control* 5 (1962) 13–24.
- [6] B. Varn, Optimal variable length codes -Arbitrary symbol cost and equal code word probability, *Information Control* (19) (1971) 289–301.
- [7] D. Altenkamp, K. Mehlhorn, Codes: Unequal probabilities, unequal letter costs, *Journal of the Association for Computing Machinery* 27 (3) (1980) 412–427.

Table 5: influence of penalty on bit minimization

Data sets	Cost	Bits	$\lambda(\%)$
Mycobacterium smegmatis	70760174	37809636	4
Amycolatopsis benzoatilytica	92010490	49654132	3
Mycobacterium rhodesiae	69421783	37563017	3
Streptomyces bottropensis	96294638	50330594	5
Mycobacterium smegmatis. MC2	74855668	40380464	5
Mycobacterium smegmatis MKD8	76738330	40782100	4
Bradyrhizobium WSM471	84667949	45961149	3
Amycolatopsis thermoflava	92416243	48336727	5
Bacillus thuringiensis serovar thuringiensis	66145783	35330943	4
Bacillus thuringiensis Bt407	68751359	37306497	3
Pseudomonas aeruginosa 9BR	72737300	40264790	2
Bacillus thuringiensis serovar berliner	67981988	36770118	3
Bacillus thuringiensis serovar pakistani	65896779	35069999	4
Pseudomonas aeruginosa LES400	71762305	38210997	4

- [8] Y. Perl, M. R. Garey, S. Even, Efficient generation of optimal prefix code: Equiprobable words using unequal cost letters, Journal of the ACM (JACM) 22 (2) (1975) 202–214.
- [9] S. Kabir, T. Azad, A. S. M. A. Alam, M. Kaykobad, Effects of unequal bit costs on classical huffman codes, in: 17th International Conference on Computer and Information Technology, 2014, pp. 96–101.
- [10] D. Howe, M. Costanzo, P. Fey, T. Gojobori, L. Hannick, W. Hide, D. P. Hill, R. Kania, M. Schaeffer, S. S. Pierre, S. Twigger, O. White, S. Y. Rhee, Big data: The future of biocuration, Nature 455 (2008) 47–50.
- [11] F. Valentin, S. Squizzato, M. Goujon, H. McWilliam, J. Paern, R. Lopez, Fast and efficient searching of biological data resources using eb-eye, Briefings in bioinformatics 11 (4) (2010) 375–384.
- [12] C. Scott, L. Yiming, L. Chen, X. Xiaohui, Human genomes as email attachments, Bioinformatics 25 (2) (2009) 274–275. doi:10.1093/bioinformatics/btn582.

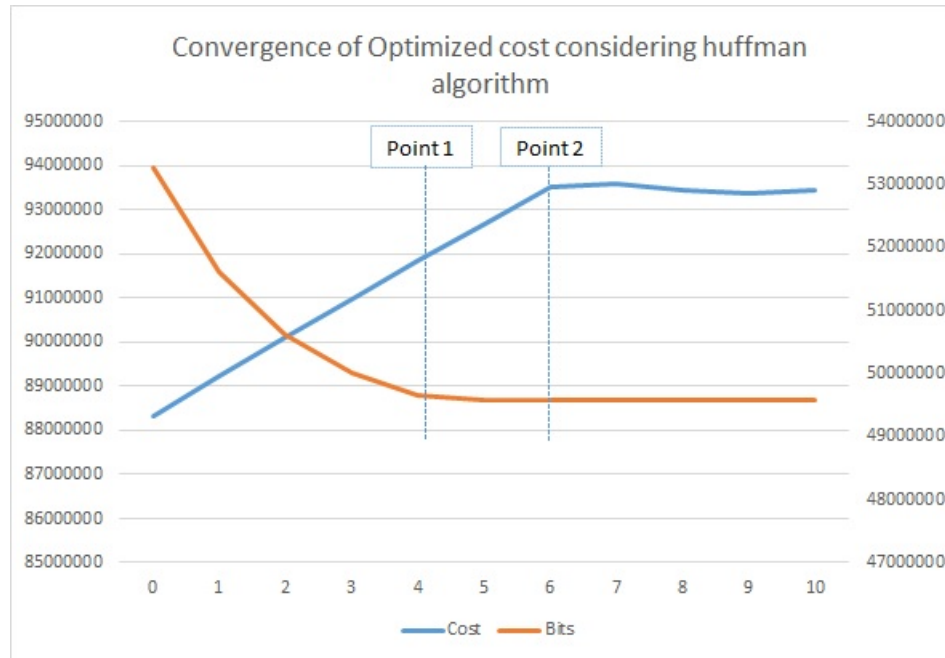


Figure 4: Convergence of Optimised cost considering Huffman algorithm

- [13] C. Ramu, S. Hideaki, K. Tadashi, L. Rodrigo, G. T. J, H. D. G, T. J. D, Multiple sequence alignment with the clustal series of programs, *Nucleic acids research* 31 (13) (2003) 3497–3500.
- [14] C. Tzu-Hao, W. Shih-Lin, W. Wei-Jen, H. Jorng-Tzong, C. Cheng-Wei, A novel approach for discovering condition-specific correlations of gene expressions within biological pathways by using cloud computing technology, *BioMed research international* 2014.
- [15] B. M. C, W. D. C, B. Pierre, Data structures and compression algorithms for genomic sequence data, *Bioinformatics* 25 (14) (2009) 1731–1738.
- [16] W. Congmao, Z. Dabing, A novel compression tool for efficient storage of genome resequencing data, *Nucleic acids research* 39 (7) (2011) e45–e45.
- [17] J. Amsterdam, Data compression with huffman coding, *BYTE* 11 (5) (1986) 98–108.

- [18] M. A. Mannan, M. Kaykobad, Block huffman coding, *Computers and Mathematics with Applications*.
- [19] M. Golin, N. Young, Prefix codes: Equiprobable words, unequal letter costs, *SIAM JOURNAL ON COMPUTING* 25 (6) (1996) 1281–1292.
- [20] M. J. Golin, C. Kenyon, N. E. Young, Huffman coding with unequal letter costs, in: *ACM Symposium on Theory of Computing*, 2002, pp. 785–791.
- [21] P. Bradford, M. Golin, L. Larmore, W. Rytter, Optimal prefix-free codes for unequal letter costs: Dynamic programming with the Monge property, *JOURNAL OF ALGORITHMS* 42 (2) (2002) 277–303.
- [22] M. J. Golin, C. Mathieu, N. E. Young, Huffman Coding with Letter Costs: A Linear-Time Approximation Scheme, *SIAM JOURNAL ON COMPUTING* 41 (3) (2012) 684–713.
- [23] M. J. Golin, G. Rote, A dynamic programming algorithm for constructing optimal prefix-free codes with unequal letter costs, *IEEE Transactions on Information Theory* 44 (5) (1998) 1770–1781.
- [24] M. Mitchell, *An introduction to genetic algorithms*, MIT press, 1998.
- [25] T. Blickle, L. Thiele, *A comparison of selection schemes used in genetic algorithms* (1995).
- [26] E. Osaba, R. Carballedo, F. Diaz, E. Onieva, I. de la Iglesia, A. Perallos, Crossover versus mutation: A comparative analysis of the evolutionary strategy of genetic algorithms applied to combinatorial optimization problems, *The Scientific World Journal* 2014.