

Reducing Power Consumption in Huffman Coding

Abstract

Recently the power consumption identified as one of the most pressing challenges of data transmission. Though a considerable reduction of size of data is obtained by Huffman code and several methods have been proposed to improve the data compression rate. These methods improves only the total number of bits without consideration of transmission cost. In this paper, we propose a new approach to reduce the power consumption, by applying genetic algorithm for minimising the total number of switches, either inter- or intra-switch. The approach starts its operation by generating a set of trees considered as initial population for the genetic algorithm, each one of these trees represent a set of codewords (solution), after that the genetic operators such as selection, crossover and mutation will be applied to the initial population in order to evolve the quality of the solutions. The amount of power consumption is proportional to the total number of switches on the whole sequence. The performance of the approach is evaluated by applying it to compress some standard biological datasets. The experiments yield that the proposed approach improves the power consumption rate considerably comparing with classic Huffman code.

Keywords:

1. Introduction

In the recent years, application of battery-powered portable devices, e.g., laptop computers, personal digital assistant (PDA), and mobile phones has increased significantly. The reliability and performance of such devices are primarily dependent on their power consumption, i.e., effective battery life. Bigger battery size could help to improve power efficiency, however, in a portable device the size of the battery is restricted by the size and weight of the device itself. The cost of providing power to both portable and non-portable devices has resulted in significant interest in power reduction. CMOS technologies were developed in order to reduce the power consumption

both in data processing and transmission. In [1], an approach was proposed to minimise power consumption in CMOS circuits. The method considers optimising the technology used to implement the digital circuits, circuits style and topology, the architecture for implementing the circuits and the algorithms that are being implemented in the devices.

As data transmission is a common operation in both portable and non-portable devices, therefore, reducing power dissipation in transmission systems is important. Authors in [2, 3] have investigated different criteria for low power data transmission. Power consumption for data transmission in IEEE 802.11 multi-hop networks and 3G mobile wireless networks were analysed in [4] and [5] respectively. In order to increase transmission speed and reduce power dissipation, 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 and obviously a significant amount power would be dissipated in the process. However, data compression techniques are widely used to reduce size of data before transmitting over serial medium to reduce transmission time and cost.

Huffman code [6] is an entropy encoding algorithm widely used for lossless data compression. For any given set of symbols and associated occurrence probabilities, Huffman algorithm generates a binary tree (Huffman tree) to encode a message with an optimal number bits. The codeword for a distinct symbol is generated by traversing the tree from the root to the leaf node representing the symbol. In this traversal process, a move towards a left and a right node is represented by 0 and 1 respectively. A major drawback of the classical Huffman code is that the whole stream must be read prior to encoding. To transmit a binary string, power dissipation is proportional to the number of switching, i.e., transition from 0 to 1 and vice versa, present in a transmitted string. Although a number of approaches like [7, 8] have extended the classical Huffman code by treating binary bits differently (with unequal letter cost), a little effort has been made to optimise the switching activities in the Huffman code to develop a power efficient Huffman code. To the knowledge of the authors, the only effort to generate Huffman code considering switching activity is seen in [9]. This approach follows a greedy strategy and works in two steps to reduce total number of switching in the

Huffman code. In the first step the switching activity between each pair of symbol (inter-switch) is reduced and then in the second step switching activity in the codeword for each symbol (intra-switch) is reduced. As the optimisation is done independently in two steps therefore it is highly likely that the gain obtained in the first step may be compromised in the second step and vice versa.

In this paper, a genetic algorithm based approach is proposed to reduce the total number of switching where inter-switch and intra-switch is optimised in a single step, i.e, a balance is made between the two different types of switching. The approach generates a set of trees, each one represents a set of codewords (solution). After that a set of genetic operators such as selection, crossover and mutation is used to evolve (create new trees) the population to minimise number of total switches. The evolution process is controlled by the total number of inter- and intra-switches. Inter-Switch between two adjacent symbols depend on the their codewords, i.e., how the codeword of the preceding symbol finished (with 0 or 1) and how the codeword for the descendent symbol started. If the end bit of the preceding symbol and the starting bit of the successor symbol differs then there exists a switch and the total number of switch is the number of times the two symbols in consideration occurs one after another. To obtain the total number Inter-Switch between pair of symbols, we used a descendent matrix which is a two dimensional matrix represents the number of occurrence of each symbol after another. This matrix is formed when the message to be transmitted is scanned to obtain the frequency of distinct symbols. The number of intra-switch for a symbol is calculated by multiplying the frequency of that symbol by the number of logic level transition (0 to 1 or 1 to 0) in the codeword of the symbol. The efficiency of the approach is evaluated by applying it to a set of biological datasets. The experiments yield that the proposed approach improves the total number of switches by $X\%$ in the best case, by $Y\%$ in the average case and by $Z\%$ in the worst case.

The rest of the paper is organised as follows: Section 2 presents the background study of the relationship of switching activity with power consumption, the relation of switching with Huffman code and the issues in biological data transmission. 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. Huffman Coding

2.2. Power consumption

2.3. Issues in Biological data transmission

The size of biological data base increase with an expanding rate and it doesn't stop for increasing. In the last decade the use of biological data history increase which evolve a high request to data from this high volume biological data base [10]. The exponential growth of these databases become a big problem to all biological data processing methods [10]. Data compression methods is the most used methods to evolve the reduce the cost of treatment of high volume data bases. The main objective of data compression methods is minimising the number of bits in the data representation. Authors in [11] proposed a new approach to genomic data compression based in suffix and common substring and code the difference with Huffman code. In [12], authors propose a general scheme for coding only the difference between the target genome and the referential one, the coding method used is Huffman. Recent work on data compression for biological data have been proposed to deal with transmission of genomic data as an email attachment [13]. Cost of data transmission is based on two point, firstly the size of the data which has been improved considerably in the recent year with the improvement of Huffman code. Secondly, is the number of switches on the generated codes. The power consumption increase linear with the number of switches.

3. GA-SO : Genetic algorithm for Switches Optimising

Genetic algorithm (GA) is one the of the most popular bio-inspired meta-heuristic algorithm inspired from the natural evolution of species [14]. It is a population based algorithm starts with the generation of a random initial population. The population contain a set of feasible solutions called individuals, that are usually far from the optimal. The GA optimisation process uses a set of natural genetic operators such as selection, crossover and mutation to converge to the optimal solution.

3.1. Preparing data and population generation

A genome is an expansion of DNA sequence which is a set of nucleotides that encodes a protein. These nucleotide are bounded together and each order define a specific protein. The transformation of DNA sequence to protein

is made by the cellular machinery, each three nucleotides are interpreted to a specific amino acid this three nucleotide are called triplets [15]. Huffman algorithm read the whole message to compute frequencies, the proposed algorithm read the whole genome and cut it into triplet to generate the frequencies of the triplet and by the way to generate the adjacent matrix which represents the frequencies precedence between triplet. The initial population is generated randomly, each solution on the population represents a tree that generates a set of codes equal to the total triplets number.

3.2. Selection

The first operator of the genetic algorithm is the selection. The main objective of the selection is to choose the part of the current population to be a candidate for the different genetic operators in order to breed the next generation population. Many selection techniques have been proposed in the literature [16]. In this approach the process of natural selection is maintained. First we generate a random pair numbers *rand* between 2 and the population size, this number represent how many parents will be processed. After that, an unbiased random selection of *rand* individuals (solutions) from the population is made.

3.3. Crossover

The main operator of the GA is the crossover, allows to construct new solutions from the selected part of the population [17]. The selected solutions are ranked by fitness and crossover two by two from high fitness to low. The main objective of the crossover is to benefits from the two good solutions in order to generate a better solution. An internal node for each tree (solution) is selected (*node1,node2*), these two nodes must have the same number of leaf nodes (contain the same number of codes). The crossover operation will create two new trees (solutions), the first child (second) contain the nodes that are not children of the *node1* (*node2*) and we replace the children of *node1* (*node2*) by the children of *node2* (*node1*) (see fig.1).

3.4. Mutation

The mutation operator changes the positions of two leaf nodes of the generated children (result of the crossover) [18]; this introduce the diversity in the search process, this diversification strategy allow the algorithm do conference to the global optimum. The algorithm for this operator goes through the leaf nodes and change the position of two leaf nodes from a parent

to another and selects the ones to change according to a fixed mutation rate (see fig.1).

3.5. Population update

The crossover operation aims to generate new solutions from the current solution to built the second generation of the population. Furthermore, the mutation provide the diversity on the search space solution by changing the position of nodes on the same tree. After these two operations, the next step of the genetic algorithm is two to breed the new population (see fig.2). Firstly the new solutions (children) are added to the current population, after that the population is ranked by fitness. The algorithm remove the worst solutions until the initial size of the population is achieved. The algorithm genetic repeat these operators until the stopped criteria is achieved, in this algorithm the stopped criteria is when the objective function (number of switches) stop decreasing.

Algorithm 1 Switches optimising Huffman codes

Require: Textual representation of a Genome

Ensure: Low switches Huffman codes

- 1: Generate triplet frequencies and adjacent matrix
 - 2: Generate the initial population of trees
 - 3: **repeat**
 - 4: Select part of the population
 - 5: Generate children candidates via crossover
 - 6: Mutate children
 - 7: Add children to population
 - 8: Rank by fitness
 - 9: Remove worst solutions until population limit
 - 10: **until** The switches number stop decreasing
-

4. Experiments and Comparison

The approach has been evaluated with different real genomic biological data. Different genomes size have been used for the evaluation, from small genome to human chromosome in order to evaluate the efficiency of the approach. The datasets were downloaded as FASTA text file from a recent version of The National Center for Biotechnology Information (NCBI) available

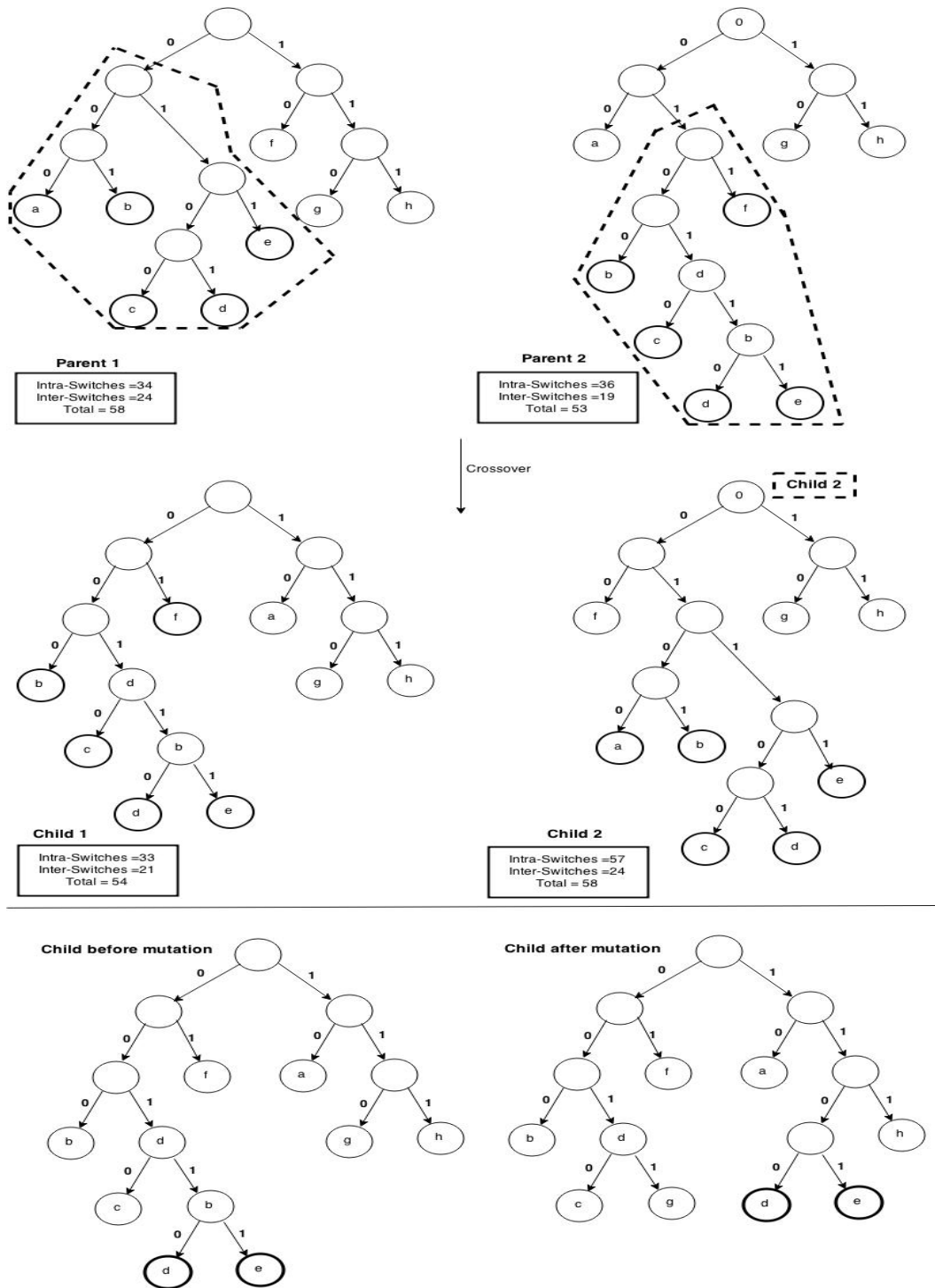


Figure 1: The Crossover operator

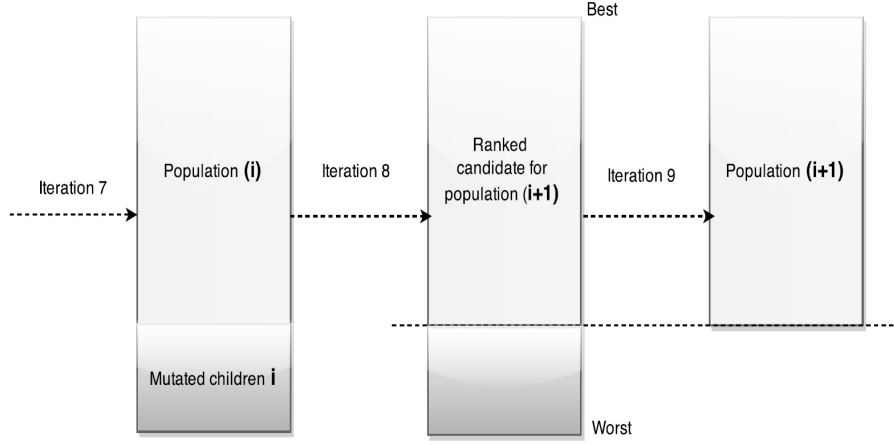


Figure 2: Population update for genetic algorithm

on (<http://www.ncbi.nlm.nih.gov>) [19]. The used genomes are described in table 1, the size are given in Megabyte and the reference on the NCBI database.

Table 2 presents the results obtained by the proposed approach and classical Huffman code. The results are presented with the total number of switches and the improvement rate. The results show that the proposed approach optimises the power consumption of the whole genome by reducing the total number of inter- and intra-switches. In the best case the proposed approach improves power consumption over classical Huffman code by x , in the worst case by x , and on an average by x . The proposed approach optimises the total number of switches by swapping finding the best arrangement of nodes on the codewords tree with the best allocation of these codewords to the different frequencies (triplet). The approach continues swapping between the population individuals until a balance is found between inter-switches and intra-switches.

5. Conclusion

In this paper we described how the genetic algorithm can be used to improve the power consumption in biological data transmission. The advantage of the proposed approach is that it minimises considerably the number of switches which affect directly the power consumption for data transmission.

Table 1: Datasets description

Data sets	Name	Size (MB)	Reference
Genome 1	Mycobacterium smegmatis	6.66	CP009496
Genome 2	Amycolatopsis benzoatilytica	8.30	NZ_KB912942
Genome 3	Mycobacterium rhodesiae NBB3	6.11	CP003169
Genome 4	Streptomyces bottropensis ATCC 25435	8.54	NZ_KB911581
Genome 5	Mycobacterium smegmatis str. MC2 155	6.66	CP009494
Genome 6	Mycobacterium smegmatis MKD8	6.76	NZ_KI421511
Genome 7	Bradyrhizobium WSM471	7.42	NZ_CM001442
Genome 8	Amycolatopsis thermoflava N1165	8.27	NZ_CM001442
Genome 9	Bacillus thuringiensis Bt407	5.74	NZ_CM000747
Genome 10	Bacillus thuringiensis serovar thuringiensis	6.03	NZ_CM000748
Genome 11	Pseudomonas aeruginosa 9BR	6.48	NZ_AFXI01000001
Genome 12	Bacillus thuringiensis serovar berliner ATCC	5.97	NZ_CM000753
Genome 13	Bacillus thuringiensis serovar pakistani	5.75	NZ_CM000750
Genome 14	Pseudomonas aeruginosa LES400	6.28	CP006982
Genome 15	Mus musculus chromosome 1	25.58	GL456087
Genome 16	Danio rerio chromosome 1	56.14	CM002885
Genome 17	Homo sapiens chromosome 18	76.64	CM000680
Genome 18	Homo sapiens chromosome 22	99.94	CM000684

Table 2: Comparison of performance among classical Huffman code, CCA, and OCCA without penalty

	Huffman Algorithm	P	
Genome 1	18.16	10.05	44.65
Genome 2	24.66	15.63	36.61
Genome 3	18.46	10.66	42.25
Genome 4	25.18	16.26	35.42
Genome 5	19.24	16.08	16.42
Genome 6	19.61	14.96	23.71
Genome 7	22.40	13.21	41.02
Genome 8	23.87	17.70	25.84
Genome 9	17.66	10.04	43.14
Genome 10	18.14	9.89	45.47
Genome 11	19.63	11.48	41.51
Genome 12	17.59	11.86	32.57
Genome 13	17.39	10.93	37.14
Genome 14	18.53	11.64	37.18
Genome 15	78.19	50.60	35.28
Genome 16	168.98		
Genome 17			
Genome 18			

In this paper, we identified genetic algorithm as a potential method to optimise the total number of inter- and intra-switches. The Genetic algorithm start by generating a set of feasible solutions (trees), each one contain a set of codewords. After that, the approach uses the natural genetic operators to improve the quality if the solutions by generating new solution. the evolution of the population is controlled by the number of switches both inter- and intra-switches. The proposed approach is applied to biological data sets with different sizes from small genome to a human chromosome. The approach has been compared with the classic Huffman code. The results show that the proposed approach improves the power consumption rate considerably. The experiment shows that the proposed approach improves Huffman code in the best case with *best*, *average* in the average case, and *worst* case. For the future we hope to improve the allocation of codewords to different frequencies

in orderer to improve the total number of switches.

References

- [1] A. P. Chandrakasan, R. W. Brodersen, Minimizing power consumption in digital cmos circuits, *Proceedings of the IEEE* 83 (4) (1995) 498–523.
- [2] S. Gregori, Y. Li, H. Li, J. Liu, F. Maloberti, 2.45 GHz power and data transmission for a low-power autonomous sensors platform, in: *Proceedings of the 2004 International Symposium on Low Power Electronics and Design (ISLPED '04)*, 2004, pp. 269–273.
- [3] D. Yates, A. Holmes, A. Burdett, Optimal transmission frequency for ultralow-power short-range radio links, *IEEE Transactions on Circuits and Systems* 51 (7) (2004) 1405–1413. doi:10.1109/TCSI.2004.830696.
- [4] L. Wang, A. Ukhanova, E. Belyaev, Power consumption analysis of constant bit rate data transmission over 3g mobile wireless networks, in: *11th International Conference on ITS Telecommunications (ITST)*, 2011, pp. 217–223. doi:10.1109/ITST.2011.6060056.
- [5] W. Toorisaka, G. Hasegawa, M. Murata, Power consumption analysis of data transmission in iee 802.11 multi-hop networks, in: *The Eighth International Conference on Networking and Services (ICNS)*, 2012, pp. 75–80.
- [6] D. Huffman, A method for the construction of minimum-redundancy codes, *Proceedings of the Institute of Radio Engineers* 40 (9) (1952) 1098–1101. doi:10.1109/JRPROC.1952.273898.
- [7] 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.
- [8] 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.
- [9] C.-Y. Chen, Y.-T. Pai, S.-J. Ruan, Low power huffman coding for high performance data transmission, in: *International Conference on Hybrid Information Technology (ICHIT)*, 2006, pp. 71–77. doi:10.1109/ICHIT.2006.253467.

- [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] M. C. Brandon, D. C. Wallace, P. Baldi, Data structures and compression algorithms for genomic sequence data, *Bioinformatics* 25 (14) (2009) 1731–1738.
- [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.
- [13] C. Wang, D. Zhang, A novel compression tool for efficient storage of genome resequencing data, *Nucleic acids research* 39 (7) (2011) e45–e45.
- [14] M. Mitchell, An introduction to genetic algorithms, MIT press, 1998.
- [15] H. Lodish, A. Berk, S. L. Zipursky, P. Matsudaira, D. Baltimore, J. Darnell, *Molecular Cell Biology*, 4th Edition, W.H. Freeman and Co Ltd, 1999.
- [16] T. Blickle, L. Thiele, A comparison of selection schemes used in genetic algorithms (1995).
- [17] 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.
- [18] H. Mühlenbein, How genetic algorithms really work: Mutation and hill-climbing., in: *PPSN*, Vol. 92, 1992, pp. 15–25.
- [19] K. D. Pruitt, T. Tatusova, W. Klimke, D. R. Maglott, NCBI reference sequences: current status, policy and new initiatives, *Nucleic acids research* 37 (suppl 1) (2009) D32–D36.