

**A HYBRID ALGORITHM FOR
DATA COMPRESSION
USING
HUFFMAN & GENETIC
ALGORITHMS**

*A Thesis Submitted to the Council of the Science
College of Babylon University in Partial
Fulfilment of the Requirements*

of the Degree of M.SC.

In Computer Science

By

Azher Razak Hadi Witwit



Safar - ١٤٢٧

March - ٢٠٠٦

خوارزمية هجينة لضغط البيانات باستخدام

Huffman

و الخوارزمية الجينية

رسالة

مقدمة الى مجلس كلية العلوم - جامعة بابل

وهي جزء من متطلبات نيل درجة ماجستير في علوم الحاسبات

من

أزهر رزاق هادي وتوت



صفر-١٤٢٧

آذار-٢٠٠٦

Acknowledgments

I am indebted to **Allah** for helping me always to finish what I start and for helping me to present this work.

I would like to express my deep gratitude and appreciation to my supervisor **Dr. Tawfiq A. Al-Asadi** for the advice, support and encouragement he gave me during the course of this work. His suggestions, guidance and moral support in difficult times have been essential for completing this work.

I would like to express my deep gratitude and appreciation to my supervisor **Dr. Abbas Al-Himyari** for support and encouragement for me.

My deep appreciation goes to the Department of Computer Science, head of the department and members of the staff.

I also wish to thank my **special friends** for their cooperation ,and **all the persons** that helped me to finish this research.

Supervisor Certification

I certify that this thesis was prepared under my supervision at the Department of Computer Science/College of Science/Babylon University, by **Azher R. H. Witwit** as partial fulfilment of the requirements for the degree of Master of Science in computer science.

Signature :

Name : **Dr. Tawfiq A. Abbas Al-Asadi**

Title : **Assistant Professor**

Date : / / ٢٠٠٦

In view of the available recommendation , I forward this theses for debate by the examination committee.

Signature :

Name : **Dr. Abbas M. AL-Bakry**

Title : **Head of Department of Computer Science.**

Date : / / ٢٠٠٦.

Certification of the Examination Committee

We chairman and members of the examination committee, certify that we have studied the thesis entitled (**A HYBRID ALGORITHM FOR DATA COMPRESSION USING HUFFMAN & GENETIC ALGORITHMS**) presented by the student **Azher R.H. Witwit** and examined him in its contents and in what is related to it, and we have found it worthy to be accepted for the degree of Master of Science in computer science with **Very Good** degree.

Signature :
Name : **Dr. Nabeel H. Kaghed**
Title : **Professor**
Date : / / ٢٠٠٦.
(chairman)

Signature :
Name: **Dr. Sattar B. Sadkan**
Title : **Professor**
Date : / / ٢٠٠٦.
(member)

Signature :
Name : **Esraa H. Ali**
Title : **Assistant Professor**
Date : / / ٢٠٠٦.
(member)

Signature :
Name: **Dr. Tawfiq A. Abbas Al-Asadi**
Title : **Assistant Professor**
Date : / / ٢٠٠٦.
(supervisor)

Signature :
Name : Uoda

Title :

Date : / /

٢٠٠٦.

(Dean of college of science-Babylon university)

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ
لِنَرْفَعَ دَرَجَاتٍ مِّنْ نَّشَأٍ وَفَوْقَ كُلِّ ذِي عِلْمٍ عَلِيمٌ

صدق الله العظيم

سورة يوسف - آية (٧٦)

LIST OF CONTENTS

SUBJECTS	PAGE NO.
CHAPTER ONE : Overview	
١.١ Introduction.	١
١.٢ Basic Concepts in Data Compression.	٢
١.٢.١ Entropy Coding.	٢
١.٢.٢ Prefix Property.	٣
١.٢.٣ Variance.	٣
١.٣ Data Compression Strategies.	٤
١.٤ Literature Review.	٥
١.٥ Aim of research.	٨

۱.۶ Thesis Organization.	۹
CHAPTER TWO : Compression Techniques	
۲.۱ Introduction.	۱۰
۲.۲ Compression methods.	۱۱
۲.۲.۱ Run-Length Encoding.	۱۱
۲.۲.۲ The Golomb code.	۱۲
۲.۲.۳ Delta Encoding.	۱۳
۲.۲.۴ LZW.	۱۵
۲.۲.۵ Oring Bits.	۱۵
۲.۲.۶ Shannon method.	۱۶
۲.۲.۷ Huffman method.	۱۷
۲.۲.۷.۱ Main properties.	۱۹
۲.۲.۷.۲ Minimum Variance Huffman Codes.	۱۹
۲.۲.۸ Adaptive Huffman Coding.	۲۲
CHAPTER THREE: Genetic Algorithm.	
۳.۱ Introduction	۲۳
۳.۲ Basic Description.	۲۳
۳.۳ Chromosome.	۲۴
۳.۴ Outline of the Basic Genetic Algorithm.	۲۴
۳.۴.۱ Initialization.	۲۵
۳.۴.۲ Calculation of Fitness Function.	۲۶

3.4.3 Selection of the Best Individuals.	26
3.4.3.1 Roulette Wheel Selection	26
3.4.3.2 Rank Selection	27
3.4.3.3 Steady-state Selection	28
3.4.3.4 Tournament Selection	28
3.4.4 Crossover .	29
3.4.4.1 Traditional Crossover .	29
3.4.4.2 Permutation Crossover .	31
3.5 Mutation.	34
3.6 Crossover and Mutation Probability.	36
3.7 Inversion	37
3.8 Replacement	38
3.9 Encoding.	40
3.9.1 Binary Encoding.	40
3.9.2 Permutation Encoding.	40
3.9.3 Value Encoding.	41
3.10 Steady-state Algorithm.	42
3.10.1 Why Steady-state GA(ssGA) used?	44
3.11 Population Size	44
3.12 Elitism .	44
3.13 Stopping criteria.	45
3.14 Applications of GA.	45

CHAPTER FOUR : The Proposed compression System	PAGE NO.
ξ.1 Introduction	ξϕ
ξ.2 Creating the Individual & Initial Population	ξλ
ξ.3 Building Huffman Tree	ξϑ
ξ.ξ Evaluation	οο
ξ.ο Tournament Selection	ο1
ξ.ϕ Crossover Operation	οξ
ξ.γ Mutation	οϕ
ξ.λ Evaluate offspring and Replacement	οϕ
ξ.ϑ Stop Criteria	ογ
ξ.1ο Compressed Operation	ολ
ξ.11 Decompressed Operation	ϕο
CHAPTER FIVE : Result , Conclusions and Future Works	
ο.1 Result Analysis	ϕ2
ο.2 Conclusions	γ1
ο.3 Future Works	γ2

LIST OF TABLES

Table No.	Description	Page No.
Table (1-1)	Data compression techniques.	4
Table (2-1)	Huffman coding.	20
Table (3-1)	Allcode of Example 1.	60
Table (3-2)	Allcode of Example 2.	67
Table (3-3)	Allcode of Example 3.	68
Table(3-4)	Effect of Different File Size.	69
Table(3-5)	Effect of Different Chromosome Length.	69
Table(3-6)	Effect of Probability of Gene.	70
Table(3-7)	Results of Oring Files.	70
Table(3-8)	Results of DNA files.	70

LIST OF FIGURES

Figure No.	Description	Page No.
Figure (۲-۱)	Block Diagram of data compression.	۱۰
Figure (۲-۲)	Run Length Encoding Example.	۱۲
Figure (۲-۳)	Golomb Encoding.	۱۳
Figure (۲-۴)	Delta Encoding Example.	۱۳
Figure (۲-۵)	Delta encoded of signal.	۱۴
Figure (۲-۶)	Reconstructing process in oring bits.	۱۶
Figure (۲-۷)	Shannon coding.	۱۷
Figure (۲-۸)	Huffman tree (code ۱).	۲۱
Figure (۲-۹)	Huffman tree (code ۲).	۲۱
Figure (۳-۱)	PMX Crossover.	۳۱
Figure (۳-۲)	Order Crossover.	۳۲
Figure (۳-۳)	Cycle Crossover.	۳۴
Figure (۳-۴)	Mutation Effect.	۳۵
Figure (۳-۵)	Mutation Example.	۳۶
Figure (۳-۶)	Inversion.	۳۷

Figure (۳-۷)	Steady-state genetic Flow chart	۴۳
Figure (۴-۱)	Block diagram for proposed system.	۴۶
Figure (۴-۲)	The General Flowchart For the proposed system.	۴۷
Figure (۴-۳)	Huffman codes.	۴۹
Figure (۴-۴)	Conversion to the optimal solution.	۵۲
Figure (۴-۵)	Tournament Selection.	۵۳
Figure (۴-۶)	Binary tournament selection.	۵۳
Figure (۴-۷)	Cycle Crossover.	۵۴
Figure (۴-۸)	crossover operation.	۵۵
Figure (۴-۹)	Mutation Operation.	۵۶
Figure (۴-۱۰)	Niching concept (subpopulation).	۵۷
Figure (۴-۱۱)	Huffman Tree.	۵۸
Figure (۴-۱۲)	Decompression process	۶۱
Figure (۵-۱)	Variance of Example ۱.	۶۶
Figure (۵-۲)	Variance of Example ۲.	۶۷
Figure (۵-۳)	Variance of Example ۳.	۶۸

To the soul of my Father....

Το μψ Μοτηερ □.Το μψ Ωιφε □.

Το μψ χηιλδρεν Μυσταφα & Φατεμα □

ωιτη λοπε.

Abstract

Text compression plays an important role and it is an essential object to decrease storage size and increase the speed of data transmission through communication channels .

In this research , a hybrid compression system is introduced which depends on the genetic algorithm for finding the best Huffman tree and finally finding the best compression ratio and then applying an additional compression method on the results of the primary compression by applying the compression method named Oring bits , and also using a decompression algorithms for both Huffman and Oring bits.

The variety in the text characters leads to the variety in the Huffman trees and finally obtaining the best possible compression . In addition, the increase in the frequencies of the text has an affect on the compression rate .The variance has been taken as a fitness function.

The proposed method used one of the crossover operators that depends on the permutation that is cycle crossover (CX) because of the nature of the

problem under research and what it results of variety . The tournament selection has also been used for what it introduces of solution that converge at first slowly to the optimal solution and then the pressure (speed) increases at the end to reach to the optimal solution.

The genetic algorithms has been used where these algorithms research in the search space for the optimal solution and with a method not depending on the application area and without the need for using knowledge of the area.

In the proposed system, the initial population is generated from the individuals, building a Huffman tree for each individual in the population, computing a fitness function for each individual and then selecting randomly two individuals from the population for stepping into the lake of crossover. After that, the new individuals (offspring) are evaluated and compared with the population in order to take a decision of adding them or not.

The selection, crossover and evaluate operation are continued until it reaches the best Huffman tree that represents the best compression for the text.

الخلاصة

يلعب ضغط النصوص دوراً مهماً ويعتبر هدفاً أساسياً لتقليل حجم الخزن وزيادة سرعة تناقل البيانات عبر قنوات الاتصالات .

في هذا البحث نقدم نظام ضغط هجين يعتمد على الخوارزمية الجينية لأيجاد أفضل شجرة Huffman وبالتالي أفضل نسبة ضغط بطريقة Huffman وتطبيق ضغط إضافي على نتائج الضغط الأولي عن طريق تطبيق خوارزمية Oring bits , وأستخدام خوارزمية فك ضغط لكل من Oring bits و Huffman .

أن التنوع في رموز النص يؤدي الى تنوع أشجار Huffman وبالتالي الحصول على أفضل نسبة ضغط ممكن ، بالإضافة الى ذلك فإن زيادة التكرارات في النص لها تأثير على نسبة الضغط .

لقد جرى اعتماد التباين Variance كمقياس لدالة الصلاحية Fitness function ,

وجرى استخدام أحد عوامل التزاوج التي تعتمد على التباديل Permutation وهو (CX) cycle crossover وذلك بسبب نوع المشكلة ولما ينتجه من تنوع ، وكذلك جرى استخدام أنتقاء المجاميع Tournament Selection وذلك لما يقدمه من حلول تقترب في البداية ببطى الى الحل الأمثل ويزداد الضغط (السرعة) عند النهاية للوصول الى الحل الأمثل .

لقد جرى استخدام الخوارزميات الجينية ، حيث تبحث هذه الخوارزميات في فراغ الحلول عن الحل الأمثل وبطريقة غير معتمدة على مجال التطبيق وبدون الحاجة الى استخدام معرفة خاصة بالمجال .

في النظام المقترح يجري توليد مجتمع ابتدائي من الأفراد ويتم بناء شجرة Huffman لكل فرد في المجتمع ، ويتم احتساب دالة الصلاحية لكل أفراد المجتمع ، ثم بعد ذلك يتم اختيار فردين عشوائياً من المجتمع للدخول في بحيرة التزاوج ومن ثم يجري تقييم الأفراد الجدد (الأبناء) ومقارنتهم مع المجتمع لاتخاذ القرار بأضافتهما اليه أم لا ، تستمر عملية الأنتقاء والتزاوج والتقييم لحين الحصول على افضل شجرة Huffman والتي تمثل أفضل ضغط للنص .

١.١ Introduction

Data transmission and storage cost money. The more information being dealt with, the more it costs. In spite of this, most digital data are not stored in the most compact form. Rather, they are stored in whatever way make them easiest to use, such as: ASCII text from word processors, binary code that can be executed on a computer, individual samples from a data acquisition system, etc. Typically, these easy-to-use encoding methods require data files about twice as large as actually needed to represent the information. Data compression is the general term for the various algorithms and programs developed to

address this problem. A compression program is used to convert data from an easy-to-use format to one optimized for compactness. Likewise, a decompression program returns the information to its original form [۲۵,۲۹].

So **Data compression** can be defined as a technique which reduces the size of a file or data collection without affecting the information contained in the file. A compressed file, obviously, requires less storage space. Moreover, electronic transfer of a compressed file is faster than the transfer of an uncompressed file and a smaller file is consequently less susceptible to transmission errors [۲۵].

;;A۴

Finding the optimal way to compress data with respect to resource constraints remains one of the most challenging problems in the field of data compression.

There are two data compression strategies: lossy and lossless. In this research we used the second strategy.

A lossless data compression algorithm takes a string of symbols (typically ASCII characters or bytes) and translates it *reversibly* into another string, one that is *on the average* of shorter length. The words “on the average” are crucial; it is obvious that no reversible algorithm can make all strings shorter, there just aren’t enough short strings to be in one-to-one correspondence with longer strings. Compression algorithms are possible only when, on the input side, some strings, or some input symbols, are more common than others. These can then be encoded in fewer bits than rarer input strings or symbols, giving a net average gain.

There exist many, quite different, compression techniques, corresponding to different ways of detecting and using departures from equiprobability in input

strings. In this research the Huffman method has been used. The idea behind Huffman coding is simply to use shorter bit patterns for more common characters. We obtain different Huffman codes for the characters but sometimes we do not want the code word lengths to vary widely from one symbol to another (example fixed rate channels).

We can make this idea quantitative by considering the concept of *entropy*. Suppose the input alphabet has Nch characters, and that these occur in the input string with respective probabilities $p_i, i = 1, \dots, Nch$, so that $\sum p_i = 1$. And suppose that **avg** represents the average size which is equal to $\sum p_i * \text{codeword length}_i$, then the variance is equal to $\sum p_i * ((\text{codeword length}_i) - \text{avg})^2$, where the variance of a code measures how much the sizes of the individual codes deviate from the average size.

1.2 Basic Concepts in Data Compression

In this section the basic concepts of data compression are shown:-

1.2.1 Entropy Coding

We can define the entropy of a signal symbol a_i as $-P_i \log_2 P_i$. This is the smallest number of bits needed, on the average, to represent the symbol.

The amount of information contained in one, base-n symbol is:
This quantity is called the entropy of the data being transited.

$$H = - \sum_{i=1}^{n-1} p_i \log_2 p_i \quad \dots(1.1)$$

The entropy of the data depends on the individual probabilities P_i , and is smaller when all n probabilities are equal [26,28].

1.2.2 Prefix Property

The prefix code can be defined as no codeword is a prefix to any other codeword. All prefix codes are uniquely decodable. They can be decoded

instantaneously since the end of the codeword is immediately recognizable

[^, 1^].

Example:

A: 1 0

B: 0 0

C: 0 1 1

D: 1 1 1

ABADCB -> 1 0 0 0 1 0 1 1 1 0 1 1 0 0 -> ABADCB

Prefix codes are based on the following facts:

1. In an optimal code, shorter code words should be used for more frequent symbols.

2. In an optimal code, the two less frequent symbols have code words with the same length.

1.2.3 Variance

The variance of a code is defined as a measurement of how much the size of the individual codes deviates from the average size [^].

$$\text{Variance} = \sum_{i=1}^n p_i (a_i - A)^2 \dots (1.2)$$

where,

p: probability

a_i : number of bits(length code word)

A: average size

1.3 Data Compression Strategies

Table (1-1) shows two different ways that data compression algorithms can be categorized. The methods have been classified as either **lossless** or **lossy**. A **lossless** technique means that the restored data file is identical to the original. This is absolutely necessary for many types of data, for example: executable code, word processing files, tabulated numbers, etc. You cannot afford to misplace even a single bit of this type of information. In comparison, data files that represent images and other acquired signals do not have to be kept in perfect condition for storage or transmission. All real world measurements inherently contain a certain amount of noise. If the changes made to these signals resemble a small amount of additional noise, no harm is done. Compression techniques that allow this type of degradation are called **lossy**. This distinction is important because lossy techniques are much more effective at compression than lossless methods. The higher the compression ratio, the more noise are added to the data [29, 24].

Table(1-1) Data Compression Techniques

Lossless	Lossy
Run-Length Encoding	CS&Q
Huffman	JPEG
Delta	MPEG
LZW	
GOLOMB	

1.4 Literature Review

There are many researches and literature in the area of compression, but we confide the ones related to our research. Among those literatures are:

WEE WEONG NG and others [34] presented a research on a genetic algorithm for performing lossless compression on text data. The text can be unambiguously parsed into words then can be replaced by code. The genetic algorithm here approximates a solution to the optimal dictionary construction problem. Where each chromosome represents a dictionary, each gene is of length w_i , where $1 \leq w_i \leq n$. The population consists of a set of chromosomes whose genes are randomly initialised with characters of one length. The fitness function is the coding rate (b/s), where b is the number of bits in the text encoded and s is the number of symbols in the uncoded text. Strong chromosomes are those that yield lower coding rates. The crossover includes selecting, exchanging, fusion, and diffusion. The selection strategy for choosing parent chromosome is random. The exchange operator exchanges genes between two parents chromosomes, the strategy followed here is that the weak genes of one parent chromosome are placed by strong genes of another parent chromosome. The fusion operator fuses two genes into a single gene. The diffusion operator splits genes into smaller genes. Both the fusion and diffusion strategies select weak gene for fusion and diffusion respectively. The compression algorithm used here is the Huffman algorithm.

Jessica Socha [14] presented a research the primary purpose of which is to outline and describe the effectiveness of meta-heuristics (genetic algorithms (GA), simulated annealing (SA), tabu search (TA)) for the lossless compression problems. In our research we focused on the genetic algorithm. The compression used is the Huffman algorithm. The GA was used to derive an

optimal tree for lossless compression of an alphabet of the English language. This work included a chromosome representation where each node in the tree would always contain only a single letter. The selection schema based on the roulette wheel and the elitist strategy. The fitness of a chromosome is calculated by using redundancy of the coding schema it represents. To perform the crossover of two chromosomes, a random subtree from each chromosome is selected and swapped. Once a subtree has been selected for swapping, its parents are recorded. The subtree is removed from the tree, and attached to the parent of the opposite subtree. The researcher had used a repair algorithm of the subtree to be swapped before they are attached to the opposite parent to avoid a resulting tree with too few or too many leaf nodes, and possibly duplicated characters. To mutate a chromosome two terminal nodes of a single tree are selected at random and swapped. The GA's optimal parameters were determined after running many test sets. Empirically, the GA shows the best performance with a small population size. The GA also produced better results when it was given more time to run. As a conclusion of the research as a comparative study of meta-heuristics, it was shown that the GA/Tabu hybrid is the most effective meta-heuristics for the lossless compression problem.

H. J. Martinez and others [14], presented a technique for image compression using binary Evolutionary Cellular Automata (ECA) which was evolved with a Genetic Algorithm (GA). The method is applied to black and white fingerprint images. The Genetic algorithm is guided by a fitness function consisting of a similarity measure between the configurations of the ECA and the target image. When a suitable approximation of the target image is achieved it can be codified with the rule numbers and logical operations linking the ECA. In this way the original image undergoes a compression process since it will be represented by a considerably smaller bit amount. Typical compression rates of the order of 10:1 are achieved in the experiments. A codified or compressed

image can be regenerated recovering a very good approximation with a typical similarity of 96% to the original image. The method is compared in speed and compression rate with other commercially and widely used compression algorithms.

Andrew Armstrong and Jamming Jiang [3], described a novel approach to image quantization by using a genetic algorithm to find near-optimal code words. The proposed system applies genetic operators directly on pixel blocks thus avoiding any conversion to a conventional format. To operate on the data this way, operators for fitness, combination and mutation have been tailored to produce near optimal results in terms of high visual and statistical quality. The popular Competitive Learning Neural Network was chosen to benchmark the system, because in terms of format, data and operating it is the closest stochastic technique to the proposed algorithm. The ability for the Genetic Algorithm to reach a near global optimum is demonstrated by the quality of the reconstructed images produced, when compared to the Neural Network for the same compression ratios. A key feature of the Genetic Algorithm is the ability to produce codebooks that contain sufficient words to fully convey detail, flat regions and shading. On the other hand the Competitive Learning Neural Network, when used under the same settings, produces high levels of detail at the cost of heavy artefacts in smooth and shaded areas. When both compressors are measured by PSNR values of their reconstructed image quality, the proposed algorithm constantly outperforms the benchmark for a group of test samples, which is evidenced by our extensive experimental reports.

Al-Asadi Tawfiq A.[30] proposed a compression system based on genetic algorithm (steady state) for finding optimal compression ratio. He described the design and implementation of novel compression and decompression algorithms .

The GA was applied to obtain an optimal order which gives a highest compression ratio, where the compression algorithm treats the data as blocks and the compression ratio was increased through rearranging the order of the blocks.

The proposed compression method and the compression system were examined by using different images with different sizes and they gave good compression ratios, and this depended on image nature and many parameters such as block size, chromosome length and population size .

1.5 Aim of Research

Finding the best Huffman tree (best compression) ; and we will get on Codewords with approximate lengths (Minimum Variance Huffman codes) ; and this will help in the transmission operation .

1.6 Thesis Organization

This thesis has been organized into five chapters:-

Chapter One :- Overview.

This chapter gives an introduction to data compression , and also illustrates some concepts in data compression. It gives an overview of data compression strategies and literature review.

Chapter Two:- Compression Techniques.

This chapter describes the compression techniques and precisely huffman method.

Chapter Three:- Genetic Algorithm.

This chapter describes the genetic algorithm and genetic operators.

Chapter Four:- The Proposed compression System.

This chapter describes and discusses the components of the proposed method.

Chapter Five:- Implementation , Conclusions and Future Works.

This chapter presents some experimental results using the proposed method and provides the conclusions of the research.

2.1 Introduction

Finding the optimal way to compress data with respect to resource constraints remains one of the most challenging problems in the field of source coding. Data compression (or source coding) is the process of creating binary representations of data which requires less storage space than the original data [34, 12].

Figure(2-1) shows the block diagram of data compression:-

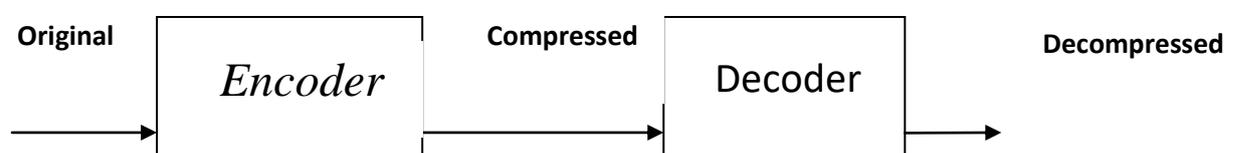


Figure (2-1) Block Diagram of Data Compression

The reason of compression is:

1- Conserving storage space.

2- Reducing time for transmission:-

Faster to encode, send, and decode them to send the original.

3- Progressive transmission:

Some compression techniques allow us to send the most important bits first so we can get a low resolution version of some data before getting the high fidelity version.

4- Reducing computation:

Use less data to achieve an approximate answer [9].

2.2 Compression Methods

There are many compression methods; we will mention some of them.

2.2.1 Run-Length Encoding

Data files frequently contain the same character repeated many times in a row. For example, text files use multiple spaces to separate sentences, indent paragraphs, format tables & charts, etc. Digitized signals can also have runs of the same value, indicating that the signal is not changing. For instance, an image of the nighttime sky would contain long runs of the character or characters representing the black background. Likewise, digitized music might have a long run of zeros

between songs. Run-length encoding is a simple method of compressing these types of files [ξ,ϒ⁹,ϑ⁸,⁸].

Figure (ϒ-ϒ) illustrates run-length encoding for a data sequence having frequent runs of zeros. Each time a zero is encountered in the input data, two values are written to the output file. The first of these values is a zero; a flag to indicate that run-length compression is beginning. The second value is the number of zeros in the run. If the average run-length is longer than two, compression will take place. On the other hand, many single zeros in the data can make the encoded file larger than the original. Many different run-length schemes have been developed. For example, the input data can be treated as individual bytes, or groups of bytes that represent something more elaborate, such as floating point numbers. Run-length encoding can be used on only one of the characters (as with the zero below), several of the characters, or all of the characters.

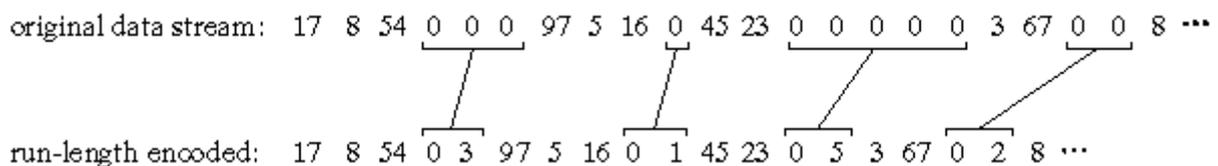


Figure (ϒ-ϒ) Run Length Encoding Example

Example of run-length encoding, Each run of zeros is replaced by two characters in the compressed file: a zero to indicate that compression is occurring, followed by the number of zeros in the run.

2.2.2 Golomb Code

The Golomb code for non-negative integer n , can be an effective Huffman code. The code depends on the choice of a parameter b . The first step is to compute the two quantities :

$$q = \left\lfloor \frac{n-1}{b} \right\rfloor, \quad r = n - qb - 1 \quad \dots(2.1)$$

Following which, the code is constructed of two parts, the first is the value of $q+1$, coded in unary, and the second, the binary value of r coded in either

bits (for the small remainders) or in $\left\lceil \frac{\log b}{2} \right\rceil$ bits (for the large ones). $\left\lfloor \frac{\log b}{2} \right\rfloor$

Choosing $b=3$, e.g. produces three possible remainders 0, 1, and 2. They are coded 0, 10, and 11, respectively. Choosing $b=5$ produces the five remainders 0 through 4, which are coded 0, 10, 100, 101, and 110 [^].

Figure (2-3) shows some examples of the Golomb code for $b=3$ and $b=5$.

n:	1	2	3	4	5	6	7	8
b=3:	0 0	0 10	0 11	10 0	10 10	10 11	110 0	110 10
b=5:	0 00	0 01	0 100	0 101	0 110	10 00	10 01	10 100

Figure(2-3) Golomb Encoding

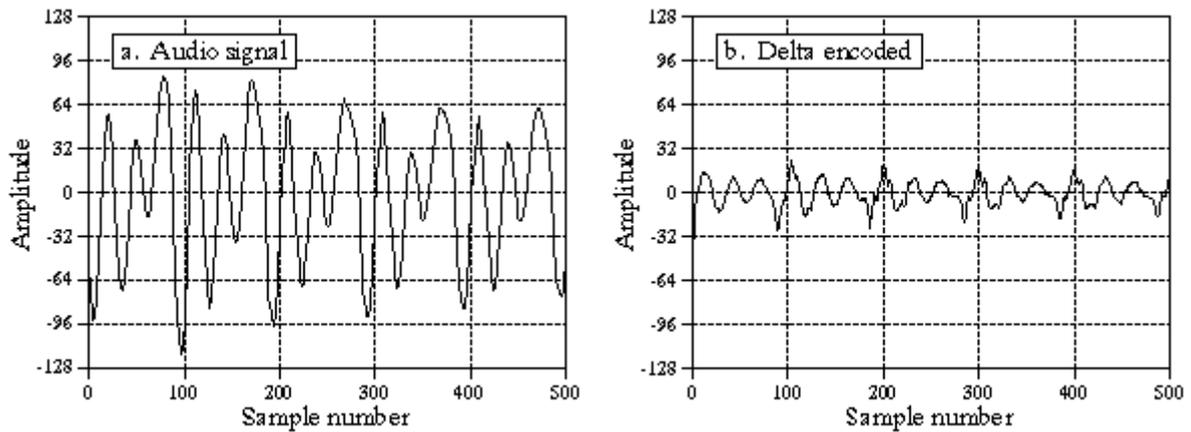
2.2.3 Delta Encoding

In science, engineering, and mathematics, the Greek letter delta is used to denote the change in a variable. The term delta encoding, refers to several techniques that store data as the difference between successive samples (or characters), rather than directly storing the samples themselves. Figure (2-4) shows an example of how this is done. The first value in the delta encoded file is the same as the first value in the original data. All the following values in the encoded file are equal to the difference (delta) between the corresponding value in the input file, and the previous value in the input file [29].

original data stream:	17	19	24	24	24	21	15	10	89	95	96	96	96	95	94	94	95	93	90	87	86	86	...	
		↓	↓	↓	↓	...																		
		move	delta	delta	delta	...																		
delta encoded:	17	2	5	0	0	-3	-6	-5	79	6	1	0	0	-1	-1	0	1	-2	-3	-3	-1	0	...	

Figure(2-4) Delta Encoding Example

Delta encoding can be used for data compression when the values in the original data are *smooth*, that is, there is typically only a small change between adjacent values. This is not the case for ASCII text and executable code; however, it is very common when the file represents a signal. For instance, Figure(2-5a) shows a segment of an audio signal, digitized to 16 bits, with each sample between -32768 and 32767. Figure(2-5b) shows the delta encoded version of this signal. The key feature is that the delta encoded signal has a lower amplitude than the original signal. In other words, delta encoding has increased the probability that each sample's value will be near zero, and decreased the probability that it will be far from zero. This uneven probability is just the thing that Huffman encoding needs to operate. If the original signal is not changing, or is changing in a straight line, delta encoding will result in runs of samples having the same value.



Figure(2.2.3) Delta Encoded of Signal

Figure (2.2.3a) is an audio signal digitized to 8 bits. Figure (2.2.3b) shows the delta encoded version of this signal. Delta encoding is useful for data compression if the signal being encoded varies slowly from sample-to-sample.

2.2.4 LZW [29, 18, 13, 14]

LZW is named after Abraham Lempel, Jakob Ziv and Terry Welch, the scientists who developed this compression algorithm. It is a lossless 'dictionary based' compression algorithm. Dictionary based algorithms scan a file for sequences of data that occur more than once. These sequences are then stored in a dictionary and within the compressed file, references are put wherever repetitive data occur.

Savings arise from the fact that pointer and length can be encoded more efficiently than the string it represents. Compression efficiency increases with the length of the text.

2.2.0 Oring Bits [A]

This method starts with a sparse string L_1 of size n_1 bits. In the first step, L_1 is divided into k substrings of equal size. In each substring all bits are logically Ored, and the results (one bit per substring) become string L_2 , which will be compressed in step 2. For example, let substrings of L_1 , that we divide into 16 substrings of size 4 each are:

$$L_1 = \dots | \dots | \dots | 1100 | \dots | \dots | \dots | 1000 | \dots$$

$$| \dots | \dots | \dots | 0100 | \dots | \dots | \dots | \dots$$

After Oring each 4-bit substring we get the 16-bit string
 $L_2 = \dots 1 | \dots 1 | \dots | 1000.$

In step 2, the same process is applied to L_2 , and the result is the 4-bit string $L_3 = 1101$, which is short enough so no more compression steps are needed. After deleting all zero substrings in L_1 and L_2 we end up with the three short strings

$$L_3 = 1101, L_2 = \dots 1 | \dots 1 | 1000, L_1 = 1100 | 1000 | 0100,$$

The output stream consists of seven 4-bit sub strings instead of the original 16 (A few more numbers are needed, to indicate how long each sub string is).

The decoder works differently (this is an asymmetric compression method). It starts with L_3 and considers each of its 4 bits a pointer to a sub string of L_2 and each of its 1 bits a pointer to a sub string of all zero that's not stored in L_2 . This way string L_1 can be reconstructed from L_3 and string L_2 , in

turn, from L_r . Figure (2-6) illustrates this process. The substrings shown in square brackets are the ones not contained in the compressed stream.

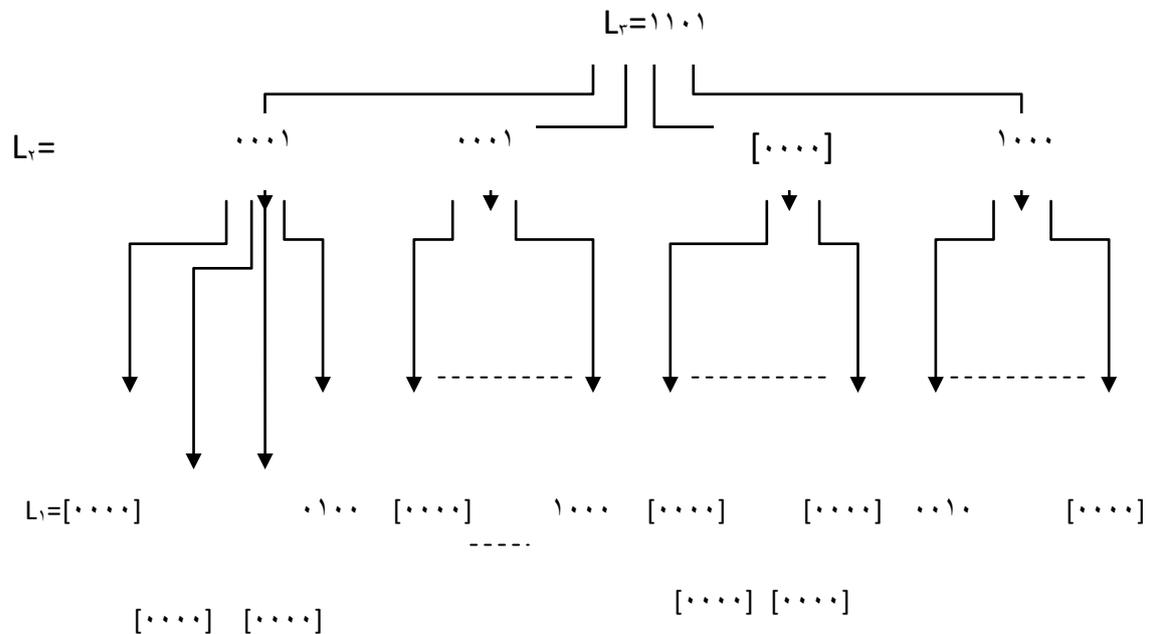


Figure (2-6): Reconstructing Process in Oring Bits.

2.2.6 Shannon Method [^]

This was the first method developed for finding good, variable size code. We start with a set of n symbols with known probabilities (or frequencies) of occurrence. The symbols are parsed arranged in descending order of their probabilities. The set of symbols is then divided into two subsets that have the same (or almost the same) probabilities. All symbols in one subset get assigned codes that start with a 0, while the codes of the symbols in the other subset start a 1. Each subset is then recursively divided into two, and the second bit of all

the code is determined in a similar way. When a subset contains just two symbols, their codes are distinguished by adding one more bit to each. The process continues until no more subsets remain.

The following example is illustrated in the figure(۲-۷):

	<u>final code</u>	<u>Steps</u>	<u>Probab.</u>
۱.	۰.۲۵	۱	۱
۲.	۰.۲۵	۱	۰
۳.	۰.۱۵	۰	۱
۴.	۰.۱۵	۰	۰
۵.	۰.۱۰	۰	۰
۶.	۰.۱۰	۰	۰
۷.	۰.۰۵	۰	۰

Figure (۲-۷) Shannon Coding

۲.۲.۷ Huffman Method [۸, ۲۹, ۱۸, ۲۶]

In ۱۹۵۱, David Huffman and his MIT information theory classmates were given the choice of a term paper or a final exam. The professor, Robert M. Fano, assigned a term paper on the problem of finding the most efficient binary code. Huffman, unable to prove any codes were the most efficient, was about to give up and start studying for the final when he hit upon the idea of using a frequency-sorted binary tree, and quickly proved this method the most efficient.

The algorithm of Huffman is:-

A fast way to create a Huffman tree is to use the heap data structure, which keeps the nodes in partially sorted order according to a predetermined criterion. In this case, the node with the lowest weight is always kept at the root of the heap for easy access.

Creating the tree:

1. Start with as many leaves as there are symbols.
2. Push all leaf nodes into the heap.
3. While there is more than one node in the heap:

-Remove two nodes with the lowest weight from the heap.

-Put the two nodes into the tree, noting their location.

-If parent links are used, set the children of any internal nodes

to point at their parents.

-Create a new internal node, using the two nodes as children

and their combined weight as the weight.

-Push the new node into the heap.

4. The remaining node is the root node.

Note it may be beneficial to generate codes with minimum length variance, since in the worst case when several long codewords have to be transmitted in a row it may lead to unwanted side effects (for example if we use a buffer and constant rate transmitter, it increases the minimum buffer size). To reduce variance every newly generated node must be favoured among same

weight nodes and placed as high as possible. This will balance the length, keeping same average rate.

2.2.7.1 *Main Properties For Huffman Method [10]*

The frequencies used can be generic ones for the application domains that are based on average experience, or they can be the actual frequencies found in the text being compressed. (This variation requires that a [frequency table](#) or other hints as to the encoding must be stored with the compressed text; implementations employ various tricks to store these tables efficiently).

Huffman coding is optimal when the probability of each input symbol is a negative power of two. Prefix-free codes tend to have slight inefficiency on small alphabets, where probabilities often fall between these optimal points. Expanding the alphabet size by coalescing multiple symbols into "words" before Huffman coding can help a bit. The worst case for Huffman coding can happen when the probability of a symbol exceeds 2^{-1} making the upper limit of inefficiency unbounded. To prevent this, [run-length encoding](#) can be used to preprocess the symbols.

2.2.7.2 Minimum Variance Huffman Codes.

The Huffman coding algorithm has some flexibility when two equal frequencies are found. The choice made in such situations will change the final code including possibly the code length of each message. Since all Huffman codes are optimal, however, it cannot change the average length [13].

For example, consider the following message probabilities, and codes in

Table (2-1):

Table (2-1) Huffman Coding

Symbol	Probability	Code 1	Code 2
A	0.2	01	10
B	0.4	1	00
C	0.2	000	11
D	0.1	0010	010
E	0.1	0011	011

$$0.4 * 1 + 0.2 * 2 + 0.2 * 3 + 0.1 * 4 + 0.1 * 4 = 2.2 \quad \text{Code 1 :}$$

$$0.4 * 2 + 0.2 * 2 + 0.2 * 2 + 0.1 * 3 + 0.1 * 3 = 2.2 \quad \text{Code 2 :}$$

Both codes produce an average of 2.2 bits per symbol, even though the lengths are quite different in the two codes. Given this choice, is there any reason to pick one code over the other?

For some applications it can be helpful to reduce the variance in the code length. The variance is defined as:

$$\text{Variance} = \sum_{i=1}^n p_i (a_i - A)^2 \dots (2.2)$$

where,

p: probability

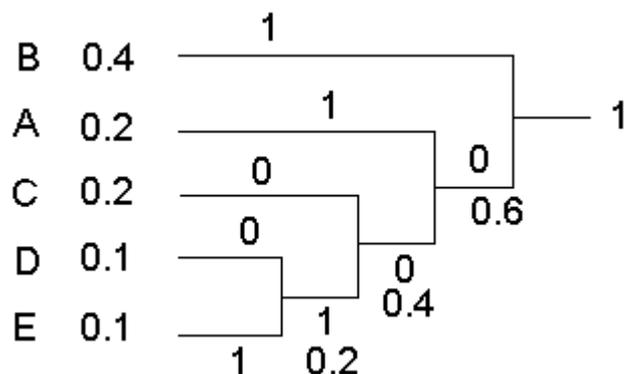
a_i: number of bits(length code word)

A: average size

$$0.4 (1-2.2)^2 + 0.2 (2-2.2)^2 + 0.2 (3-2.2)^2 + 0.1 (4-2.2)^2 + 0.1 (4-2.2)^2 = 1.36$$

$$0.4 (2-2.2)^2 + 0.2 (2-2.2)^2 + 0.2 (2-2.2)^2 + 0.1 (3-2.2)^2 + 0.1 (3-2.2)^2 = 0.16$$

With lower variance it can be easier to maintain a constant character transmission rate, or reduce the size of buffers. In the above example, code 1 clearly has a much higher variance than code 2 and it is illustrated in Figure (2-8) and Figure (2-9) respectively.



Figure(γ-λ) : Huffman Tree (code λ).

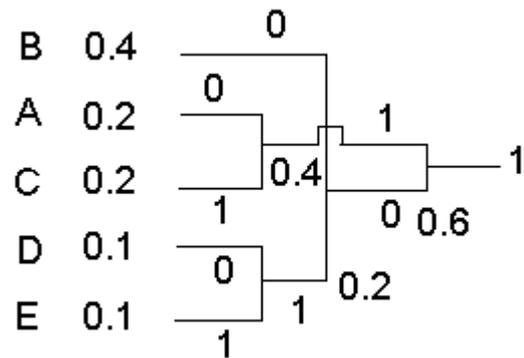


Figure (γ-α) Huffman Tree(code γ).

It turns out that a simple modification to the Huffman algorithm can be used to generate a code that has minimum variance. In particular when choosing the two nodes to merge and there is a choice based on weight, always pick the node that was created earliest in the algorithm. Leaf nodes are assumed to be created before all internal nodes. In the example above, after D and E are joined, the pair will have the same probability as C and A (0.2), but it was created afterwards, so we join C and A. Similarly we select B instead of the join AC to join with DE since it was created earlier. This will give code γ above, and the corresponding Huffman tree in Figure (γ-α).

۲.۲.۸ Adaptive Huffman Coding[۸]

Traditional Huffman Coding uses a static dictionary, which means each item in the file is encoded the same way. Typical implementations send the encoding table to the decoder as part of the encoded file.

Adaptive Huffman coding modifies the table as characters are encoded, which allows the encoder to adapt to changing conditions in the input data. Adaptive decoders do not need a copy of the table when decoding, they start with a fixed decoding table and update the table as characters are read in.

۳.۱ Introduction [۳۶,۱۰]

A **genetic algorithm (GA)** is a heuristic used to find approximate solutions for difficult to solve problems through application of the principles of evolutionary biology to computer science. Genetic algorithms use biologically-derived techniques such as inheritance, mutation, natural selection, and recombination (or crossover). Genetic algorithms are a particular class of evolutionary algorithms.

Genetic algorithms are typically implemented as a computer simulation in which a population of abstract representations (called **chromosomes**) of candidate solutions (called **individuals**) to an optimization problem evolving toward better solutions. Traditionally, solutions are represented in binary as strings of ۰'s and ۱'s, but different encodings are also possible. The evolution starts from a population of completely random individuals and happens in generations. In each generation, the fitness of the whole population is evaluated, multiple individuals are stochastically selected from the current population

(based on their fitness), modified (mutated or recombined) to form a new population, which becomes current in the next iteration of the algorithm.

२.२ Basic Description [२०]

Genetic algorithms are inspired by Darwin's theory about evolution.

Solution to a problem solved by genetic algorithms is evolved.

Algorithm is started with a **set of solutions** (represented by **chromosomes**) called **population**. Solutions from one population are taken and used to form a new population. This is motivated by a hope that the new population will be better than the old one. Solutions which are selected to form new solutions (**offspring**) are selected according to their fitness - the more suitable they are the more chances they have to reproduce.

This is repeated until some condition (for example a number of populations or improvement of the best solution) is satisfied.

२.३ Chromosome [२०]

All living organisms consist of cells. In each cell there is the same set of **chromosomes**. Chromosomes are strings of **DNA** and serve as a model for the whole organism. A chromosome consists of **genes**, blocks of DNA. Each gene encodes a particular protein. Basically it can be said, that each gene encodes a **trait**, for example color of eyes. Possible settings for a trait (e.g. blue, brown) are called **alleles**. Each gene has its own position in the chromosome. This position is called **locus**.

A complete set of genetic material (all chromosomes) is called **genome**.

A particular set of genes in genome is called **genotype**. The genotype is with later development after birth base for the organism's **phenotype**, its physical and mental characteristics, such as eye color, intelligence etc.

3.4 Outline of the Basic Genetic Algorithm [1,20]

1. **Start** Generate random population of n chromosomes (suitable solutions for the problem).
2. **Fitness** Evaluate the fitness $f(x)$ of each chromosome x in the population
3. **New population** Create a new population by repeating following steps until the new population is complete:
 - a. **Selection** Select two parent chromosomes from a population according to their fitness (the better fitness, the bigger chance to be selected).
 - b. **Crossover** With a crossover probability cross over the parents to form a new offspring (children). If no crossover was performed, an offspring is an exact copy of parents.
 - c. **Mutation** With a mutation probability mutate a new offspring at each locus (position in chromosome).
 - d. **Accepting** Place a new offspring in a new population .
4. **Replace** Use new generated population for a further run of algorithm
- o. **Test** If the end condition is satisfied, **stop**, and return the best solution in current population.

6. **Loop** Go to step 2.

As you can see, the outline of Basic GA is very general. There are many things that can be implemented differently in various problems.

The first question is how to create chromosomes, what type of encoding is chosen. This is connected with crossover and mutation which are the two basic operators of GA.

The next questions is how to select parents for crossover. This can be done in many ways, but the main idea is to select the better parents (in the hope that the better parents will produce better offspring). In this case, there is a big chance the best chromosome is lost. That is true, here is where **elitism** is used.

3.4.1 Initialization

The Algorithm is started with a set of solutions (represented by chromosomes) called population.

Each chromosome represents a possible solution to the problem. The most-used way (though not the only way) of encoding chromosomes is a binary string [1].

3.4.2 Calculation of Fitness Function [1]

An evaluation function, called fitness function, needs to be defined for a problem to be solved in order to evaluate chromosomes. The fitness of a chromosome defines how well it is suited for the problem; hence it determines its likelihood for survival from one generation to the next. A chromosome with a high fitness value is likely to be a good solution to the problem.

3.4.3 Selection of the Best Individuals:

As you already know from the GA outline, chromosomes are selected from the population to be parents to crossover. The problem is how to select

these chromosomes. According to Darwin's evolution theory the best ones should survive and create new offspring. There are many methods how to select the best chromosomes, for example roulette wheel selection, Boltzman selection, tournament selection, rank selection, steady state selection and some others [1].

Some of them will be described in this section:

3.4.3.1 Roulette Wheel Selection [20]

Parents are selected according to their fitness. The better the chromosomes are, the more chances to be selected they have. Imagine a **roulette wheel** where all chromosomes are placed in the population, every one has its big place accordingly to its fitness function.

Then a marble is thrown there and selects the chromosome. Chromosome with bigger fitness will be selected more times.

This can be simulated by the following algorithm.

1. Find the fitness f_i of each member of the population.

2. Find the sum of the fitnesses:

$$S = \sum_{i=1}^n f_i$$

3. Scale each fitness with S (convert to a probability):

$$p_i = f_i / S.$$

4. Find the cumulative probability for each individual:

$$c_i = \sum_{k=1}^i p_k \quad \text{for } i = 1, \dots, N$$

- Choose a random number r uniformly from $[0, 1]$ and select member m_i if r is such that $c_{i-1} < r \leq c_i$. (Select m_1 if $r < c_1$.) Do this N times.

3.4.3.2 Rank Selection [20]

The previous selection will have problems when the fitnesses differs very much. For example, if the best chromosome fitness is 90% of all the roulette wheel, then the other chromosomes will have very few chances to be selected.

Rank selection first ranks the population and then every chromosome receives fitness from this ranking. The worst will have fitness 1, second worst 2 etc. and the best will have fitness N (number of chromosomes in population).

After this all the chromosomes have a chance to be selected. But this method can lead to slower convergence, because the best chromosomes do not differ so much from other ones.

3.4.3.3 Steady-State Selection [20, 16]

This is not a particular method of selecting parents. The main idea of this selection is that a big part of chromosomes should survive to the next generation.

GA then works in a following way. In every generation a few chromosomes (good - with high fitness) are selected for creating a new offspring. Then some (bad - with low fitness) chromosomes are removed and the new offspring is placed in their place. The rest of population survives to new generation.

3.4.3.4 Tournament Selection [2]

This method of selection organizes the number of copies inclusively by giving the individuals with the least fitness a chance to be selected. There are many forms to perform this kind of selection. The simplest and more common is binary tournament selection that is used in (ssGA). Where two different individuals are selected randomly, the best wins and is added to the crossover pool. In the case the two fitness are equal, one of them is selected randomly. This operation is repeated for selecting the other individual that will share in the crossover operation.

1. Select k individuals (at random) from the current population.
 2. Choose the best of these individuals to be placed in the new population.
- If more than one of the selected individuals has the same fitness choose the winner at random.
3. Repeat this N times until the new population has been formed.

Notes:

- $k = 2$ is often used.
- $k = N$ would result in a new population consisting solely of the best individual(s).

The algorithm above assumes sampling with replacement [11].

3.4.4 Crossover [33]

One of the unique aspects of the work involving genetic algorithms (GAs) is the important role that crossover plays. Crossover operators produce two new chromosomes by exchanging information of the selected chromosomes. Good individuals will probably be selected several times in a generation, poor ones may not be at all. This operator is the most essential in genetic algorithms.

3.4.4.1 Traditional Crossover

After we have decided what encoding we will use, we can make a step to crossover. Crossover selects genes from parent chromosomes and creates a new offspring. The simplest way how to do this is to choose randomly some crossover point and everything before this point copy from a first parent and then everything after a crossover point copy from the second parent. In other words, Crossover selects random genes from parent chromosomes and creates a new offspring. The chromosomes of the two parents are split into two (equal or unequal) halves each. Both chromosomes are cut similarly. The halves are interchanged and combined to form the child chromosome [1,33]. The following illustrates some types of this kind of crossover.

1- One-Point Crossover (1X) [23,24]

Suppose we are using a bit-string representation and that one of the chosen pairs is 1011101011 and 11110110110 (each of length 12).

- Suppose a number at random between 1 and $k - 1$ has been chosen (here 11).
- Suppose it is \wedge . This is the crossover position (marked by the bar).

1 0 1 0 1 1 0 1 | 0 0 1 1

1 1 1 1 0 0 1 1 | 0 1 1 0

- All bits from position \wedge to the end of the string are swapped. The new strings will be (bar left in to help clarify things).

1 0 1 0 1 1 0 1 | 0 1 1 0

1 1 1 1 0 0 1 1 | 0 0 1 1

Two-Point Crossover (2X) [\wedge, \vee]

Suppose we are using a bit-string representation and that one of

the chosen pairs is 1 0 1 0 1 1 0 1 and 1 1 1 1 0 0 1 1 (each of length 8).

- Choose two numbers at random between 1 and k.
- Suppose they are \wedge and \vee . These are the crossover positions (indicated

by the bars)

1 0 1 | 0 1 1 0 1 0 | 1 1

1 1 1 | 1 0 0 1 1 0 | 1 1 0

- Swap the bits between the markers.

1 0 1 | 1 0 0 1 1 0 | 1 1

1 1 1 | 0 1 1 0 1 0 | 1 1 0

(The idea can be generalized to n-point crossover.)

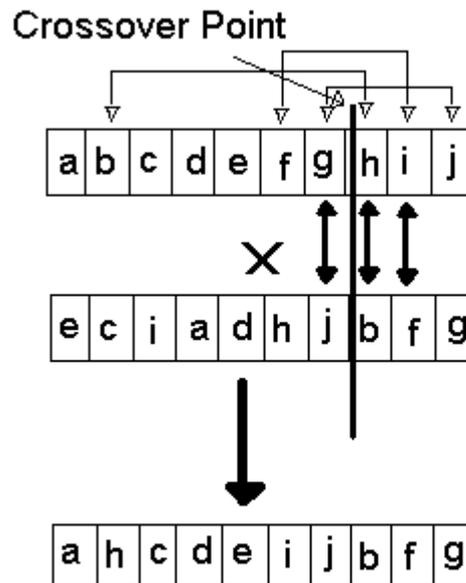
۳.۴.۴.۲ Permutation Crossover

The operators of permutation are:

a- Partially Mapped Crossover (PMX)[۲۲].

PMX is implemented as follows:

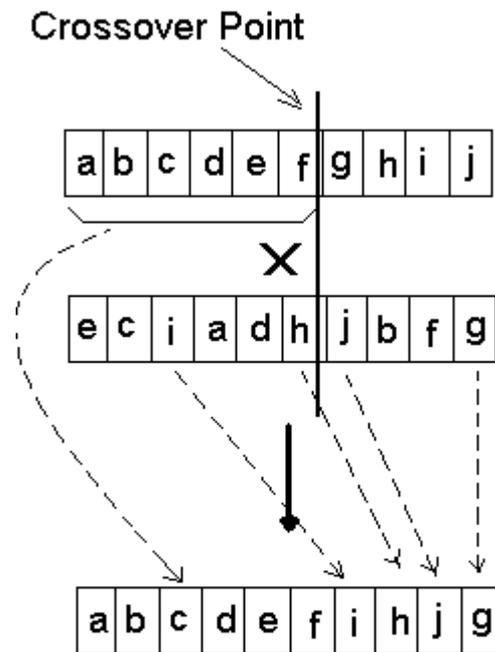
Choose a random cut point and consider the segments following the cut point in both parents as a partial mapping of the cells to be exchanged in the first parent to generate the offspring. Take corresponding cells from the segments of both parents, locate both these cells in the first parent, and exchange them. Repeat this process for all cells in the same location in the second parent will define which cells in the first parent have to be exchanged to generate the offspring. An example is shown in Figure (۳-۱), where the first parent is **abcdefghij**, the second parent is **eciadhjbfjg**, and the point cut is at the seventh cell.



b- Order Crossover (OX)[۲۲].

This kind of crossover is implemented as follows:

Choose a cut point at random. Copy the array segment to the left of the cut point from one parent to the offspring. Fill the remaining (right) portion of the offspring array by going through the second parent, from the beginning to the end and taking those elements which were left out, in order. An example is shown in Figure (۳-۲). This operator conveys a sub-placement from the first parent without any changes. And then, to resolve conflicts, compress the second parent by eliminating the cells conveyed by the first parent and shifting the rest of the cells to the left without changing their order. It then copies this compressed second parent into the remaining part of the offspring array [۲۲].



Figure(3-2) Order Crossover

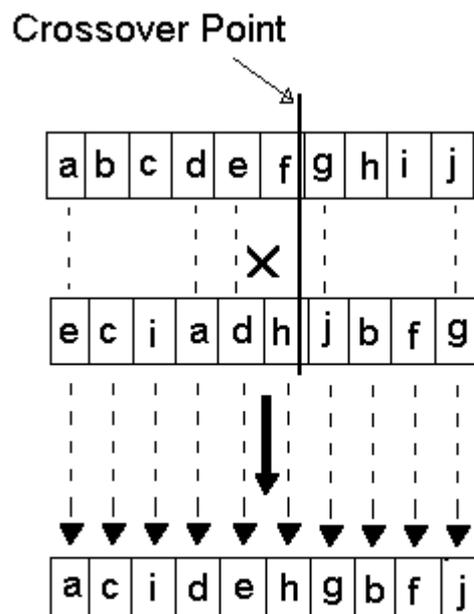
c- Cycle Crossover (CX)[۲۲].

Cycle crossover is an attempt to eliminate the cell conflicts that normally arise in crossover operators. In the offspring cells generated by cycle crossover, every cell is in the same location as one of the parents. The basic idea behind cycle crossover is to avoid cell conflicts by finding non-overlapping sets of cells to pass from the two parents. However, such non-overlapping sets of cells are not contiguous. This kind of crossover is implemented as follows:

Start with the cell in location λ (or any other reference point), and copy it to location λ of the offspring. Now consider what will happen to the cell in location λ of parent Ψ . The offspring cannot inherit this cell from parent Ψ , since location λ in the offspring has been filled. Therefore, this cell must be obtained from parent λ . Suppose this cell is located in parent λ at location X. Then it is passed to the offspring at location X. But then the cell at location X in parent Ψ cannot be passed to the offspring, so that cell is also obtained from parent λ . This process continues until we complete a cycle and reach a cell that has already been passed. Then we choose a cell from parent Ψ to pass the

offspring and go through another cycle, passing cells from parent Ψ to the offspring. Thus, in alternate cycles, the offspring inherits cells from alternate parents, and the cells are placed in the same locations as they were in the parents from which they were inherited.

An example is given in Figure (3-3). In this example, we start by passing cell **a** from parent Ψ to the offspring. Since **e** is located in the same position in parent Ψ , it cannot be passed from parent Ψ . Hence, **e** is also passed from parent Ψ . **d** is located in parent Ψ in the same position as **e** in parent Ψ . Hence **d** must also be passed from parent Ψ , and so on until we reach cell **a** again. This completes the cycle. The next cycle is taken from parent Ψ and consists of the cells **c**, **b**, **h**, **f**, and **i**. The third cycle is again from parent Ψ and consists of cells **g** and **j**.



Figure(3-3) Cycle crossover

3.5 Mutation [20,1]

After a crossover is performed, the resulting solution might fall into a local optimum. Mutation changes randomly the new offspring. (In Figure(3-4), bits are randomly toggled in case of binary strings).

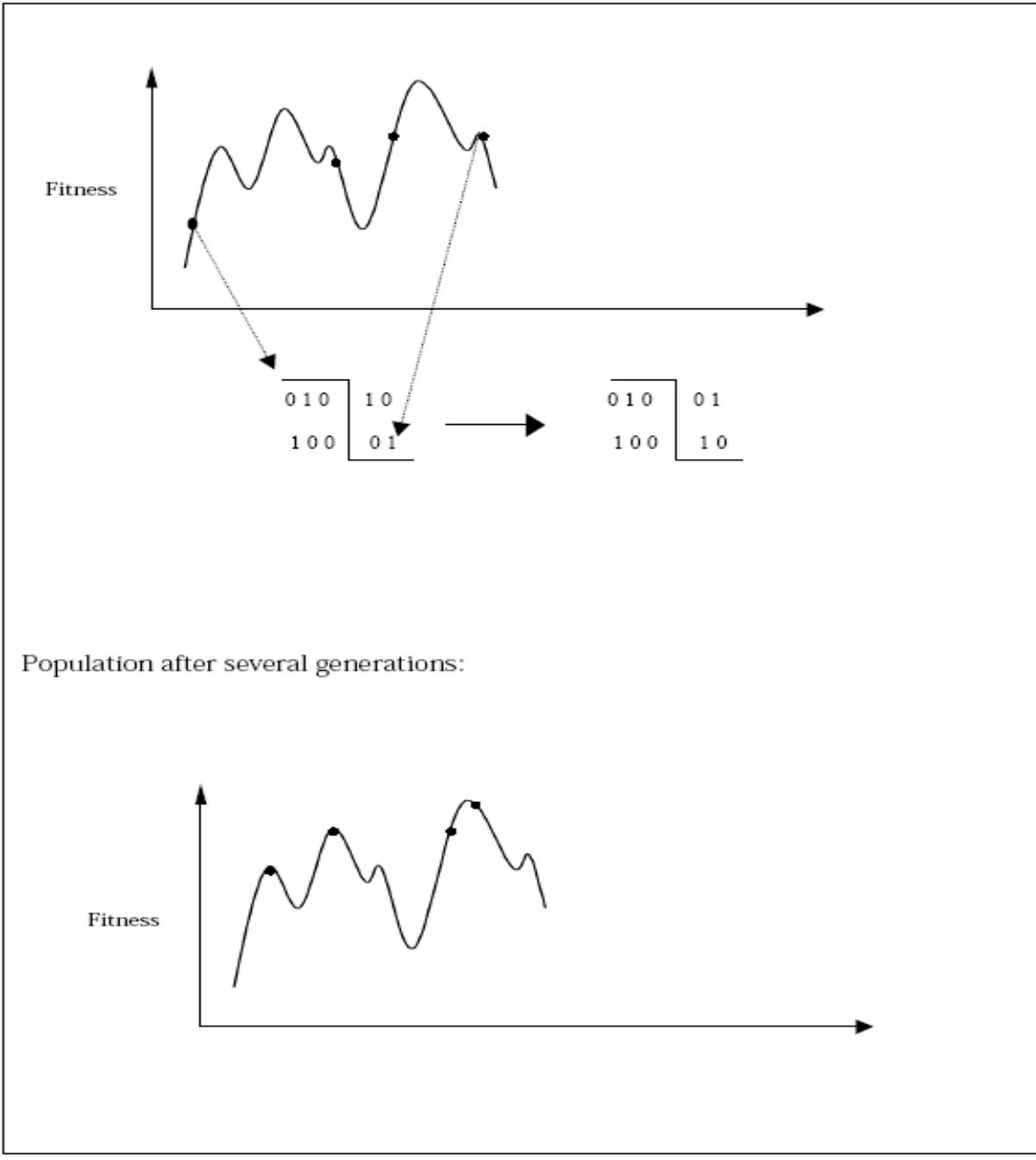


Figure (۳-۴) Mutation Effect

For binary encoding we can switch a few randomly chosen bits from 1 to 0 or from 0 to 1. Mutation can then be as follows:

Original offspring 1	110111100011110
Original offspring 2	110110100110110
Mutated offspring 1	110011100011110
Mutated offspring 2	1101101100110110

Figure (3-5) Mutation Example

The mutation depends on the encoding as well as the crossover. For example when we are encoding permutations, mutation could be exchanging two genes.

Mutation changes a gene in the chromosome with small probability P_m . Mutation is a mechanism for restoring lost and unexplored genetic material into the population and searching regions of the allele space not generated by selection and crossover.

It can be used to prevent premature convergence. This increases the diversity of genetic material.

There are many methods of mutation:

- 1- Flips the value of a gene from 1 to 0 or vice versa and as in the previous example.
- 2- Change the value of a gene.
- 3- Swap the value of two genes.

3.6 Crossover and Mutation Probability

There are two basic parameters of GA - crossover probability and mutation probability. **Crossover probability** says how often crossover will be

er	٩	٨	٧	٢	٣	٤	٥	٦	٧	٨
Inversion	٧	٧	٨	٧	٧	٨	٨	٧	٧	٨

Figure (٣-٦) Inversion

Inversion changes the sequence of genes randomly, in the hope of discovering a sequence of linked genes placed close to each other. The inversion operator greatly improves the efficiency of the crossover operator.

٣.٨ Replacement

Replacement deletes old individuals with least fitness and replaces them with new individuals of better fitness; the following are two common methods for replacement [٥]:

١- Binary tournament Replacement[٥]

This replacement method is identical to the binary tournament selection illustrated in section (٣.٤.٣.٤), where the new offspring that is evaluated from both the crossover and mutation operation is replaced with worst individual selected from the two randomly selected individuals from the population. The following mathematical form illustrates this type of replacement:

$$\text{Replace}_n = \begin{cases} \text{Ind}_i & \text{if } F(\text{ind}_i) < F(\text{ind}_j) \\ \text{Ind}_j & \text{otherwise} \end{cases} \quad \dots(٣.٧)$$

For $n=\{١,٢\}$, random number $i, j \in \{١,٢, \dots, N_{\text{pop}}\}$ $i \neq j$.

Replace_n : individual n wanted replace .

F(ind_i) : the fitness of individual_i .

F(ind_j) : fitness of individual_j .

This type of replacement methods always keeps the best individuals and by elitism size equal to one at least. To increase the probability of keeping the best individuals in the population there is another method which can be implemented as illustrated in the following replacement method.

Ψ- Triple Tournament Replacement [°]

This type is identical to the previous type of replacement except replacing the new offspring instead of the worst individual selected from three randomly selected individuals from the population. The following mathematical form illustrates that:

$$\text{Replace}_n = \begin{cases} \text{ind}_i & \text{if } F(\text{ind}_i) < F(\text{ind}_j) \text{ and } F(\text{ind}_i) < F(\text{ind}_k) \\ \text{ind}_j & \text{if } F(\text{ind}_j) < F(\text{ind}_i) \text{ and } F(\text{ind}_j) < F(\text{ind}_k) \text{ ..}(\Psi.\Psi) \\ \text{ind}_k & \text{otherwise} \end{cases}$$

For $n=\{1,2\}$, and random number $i, j, k \in \{1, 2, \dots, N_{pop}\}$ $i \neq k, j \neq k$

This type of replacement gives a higher probability from the previous type from keeping the best individuals, where it is produced implicit a elitism size equal to two or less (the best two individuals are never replaced).

3- Other Replacement Methods

In addition to the two previous methods there are other replacement methods:

- 1- Worst Replacement: Replacing worst offspring with the worst individual in the population [19].
- 2- Random Replacement: Replacing the new offspring with the individual that is randomly selected from the population [19].
- 3- In addition to that a strategy replacement may be used that takes in consider the fitness and individual's similarity to the population in order to choose the desired individual to be deleted which has a low fitness and is relatively similar to the new child [6].

3.9 Encoding [21]

Encoding of chromosomes is one of the problems, when you are starting to solve a problem with GA. Encoding deeply depends on the problem. This section will introduce some encoding, which has been used with some success.

3.9.1 Binary Encoding

Binary encoding is the most common, mainly the because first works about GA used this type of encoding.

In **binary encoding** every chromosome is a string of **bits**, 0 or 1.

The example of chromosomes with binary encoding is the following:

Chromosome A	1 0 1 1 0 0 1 0 1 1 0 0 1 0 1 1 1 0 0 1 0 1
Chromosome B	1 1 1 1 1 1 1 0 0 0 0 0 1 1 0 0 0 0 0 1 1 1 1 1

Binary encoding gives many possible chromosomes even with a small number of alleles. On the other hand, this encoding is often not natural for many problems and sometimes corrections must be made after crossover and/or mutation.

3.9.2 Permutation Encoding [21]

Permutation encoding can be used in ordering problems such as the traveling salesman problem or the task ordering problem.

In **permutation encoding**, every chromosome is a string of numbers, which *represents number in a sequence*.

Example of chromosomes with permutation encoding is:-

Chromosome A	1 0 3 2 6 4 7 9 8
Chromosome B	8 0 6 7 2 3 1 4 9

Permutation encoding is only useful for ordering problems. Even for this problems for some types of crossover and mutation corrections must be made to leave the chromosome consistent (i.e. have real sequence in it).

3.9.3 Value Encoding [21]

Direct value encoding can be used in problems, where some complicated values, such as real numbers, are used. The use of binary encoding for this type of problems would be very difficult.

In **value encoding**, every chromosome is a string of some values. Values can be anything connected to the problem, form numbers, real numbers or chars to some complicated objects.

Example of chromosomes with value encoding is:

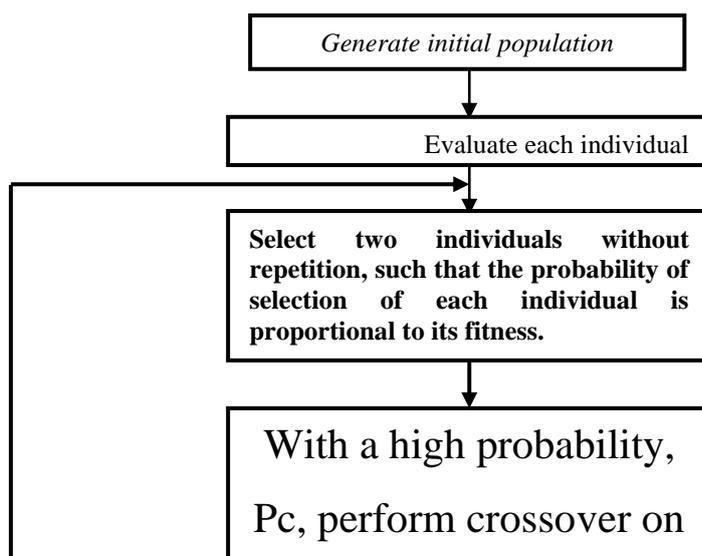
Chromosome A	1.2324 0.3243 0.4006 2.3293 2.4040
Chromosome B	ABDJEIFJDHDIERJFDLDFLFE GT
Chromosome C	(back), (back), (right), (forward), (left)

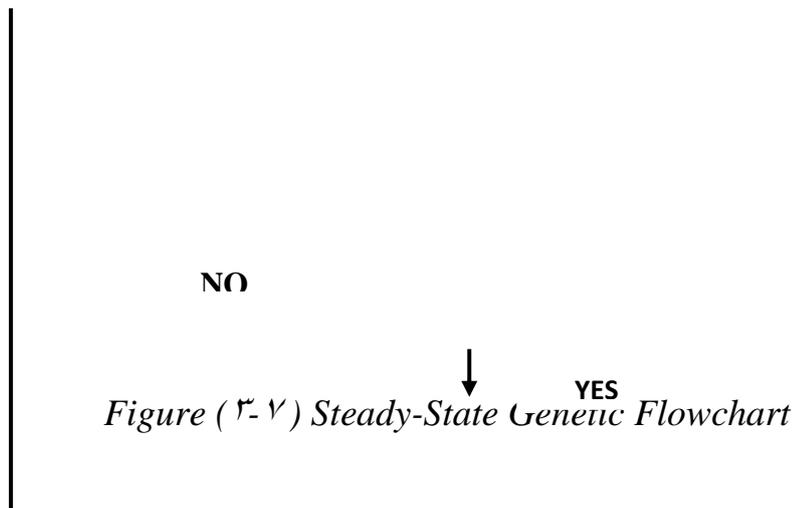
Value encoding is very good for some special problems. On the other hand, this encoding is often necessary to develop some new crossover and mutation specific for the problem.

3.10 Steady State Genetic Algorithm (ssGA).

In a GA having overlapping generations, only fraction of the individuals are replaced in each generation. The steady state algorithm is illustrated in Figure (3-9). In each generation, two different individuals are

selected as parents, based on their fitness. Crossover is performed with a high probability, P_c , to form offspring. The offspring are mutated with a low probability, P_m and inverted with probability P_I , if necessary. A duplicate check may follow, in which the offspring are rejected without any evaluation if they are duplicates of some chromosomes already in the population. The offspring that survive the duplicate check are evaluated and are introduced into the population only if they are better than the current worst member of the population, in which case the offspring replaces the worst member. This completes the generation. In the steady state GA, the generation gap is minimal, since only two offspring are produced in each generation [22, 11].





3.10.1 Why Steady-State GA (ssGA) Used ?

Simple genetic algorithm SGA replaces the entire parent population with the children; one tactic that is often added to a selection strategy is elitism. In elitism the fittest individuals pass unchanged from the parent population to the children one. This is to ensure that the fittest genes are not lost to the next generation.

Steady state GA (ssGA) replaces few individuals, while a standard GA aims at finding a single Global optimum, the application of the concept of the Niching methods (subpopulation) provide a set of solutions rather than only one solution [33].

3.11 Population Size [34]

There are also some other parameters of GA. An important parameter is also population size. **Population size** says how many chromosomes are in population (in one generation). If there are too few chromosomes, GAs have a few possibilities to perform crossover and only a small part of search space is explored. On the other hand, if there are too many chromosomes, GA slows down. Research shows that after some limit (which depends mainly on encoding and the problem) it is not useful to increase the population size, because it does not make solving the problem faster.

३.१२ **Elitism [२०]**

When creating a new population by crossover and mutation, we have a big chance, that we will loose the best chromosome. Elitism is the name of method, which first copies the best chromosome (or a few best chromosomes) to a new population. The rest is done in a classical way. Elitism can very rapidly increase performance of GA, that at least one best solution is copied without changes to a new population, so the best solution found can survive to end of run.

३.१३ **Stopping Criteria[३१]**

There are many criteria for stopping the GA process as:

- After a fixed number of generation.
- When known optimal fitness is reached.
- When GA can do no better when all individuals are identical.

३.१४ **Applications of GA**

Genetic algorithms have been used for difficult problems (such as NP-hard problems), for machine learning and also for evolving simple programs.

They have also been used for some art, for evolving pictures and music. Advantage of GAs is in their parallelism. GA is traveling in a search space with more individuals (and with genotype rather than phenotype) so they are less likely to get stuck in a local extreme like some other methods.

They are also easy to implement. Once you have some GA, you just have to write a new chromosome (just one object) to solve another problem. With the same encoding you just change the fitness function and that is all. On the other hand, choosing encoding and fitness function can be difficult. Disadvantage of GAs is in their computational time. They can be slower than some other methods. But with modern computers it is not so big problem.

Structure of proposed system

4.1 Introduction

In this chapter, Figure (4-1) illustrates the block diagram for the proposed system and its algorithms, and it has reviewed the system mechanism in compression based on some examples, where it illustrates the compression operation through building the trees and codes and applying them .

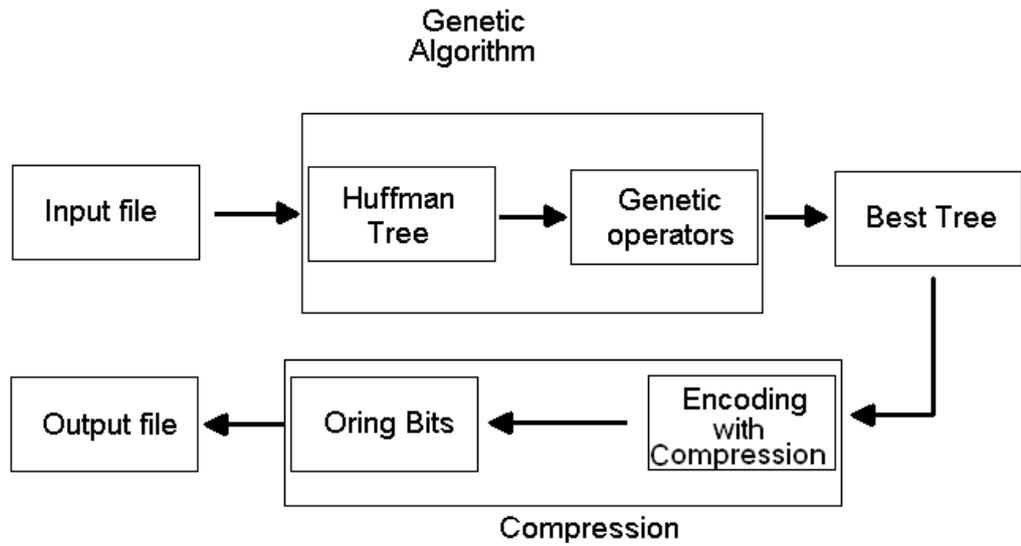
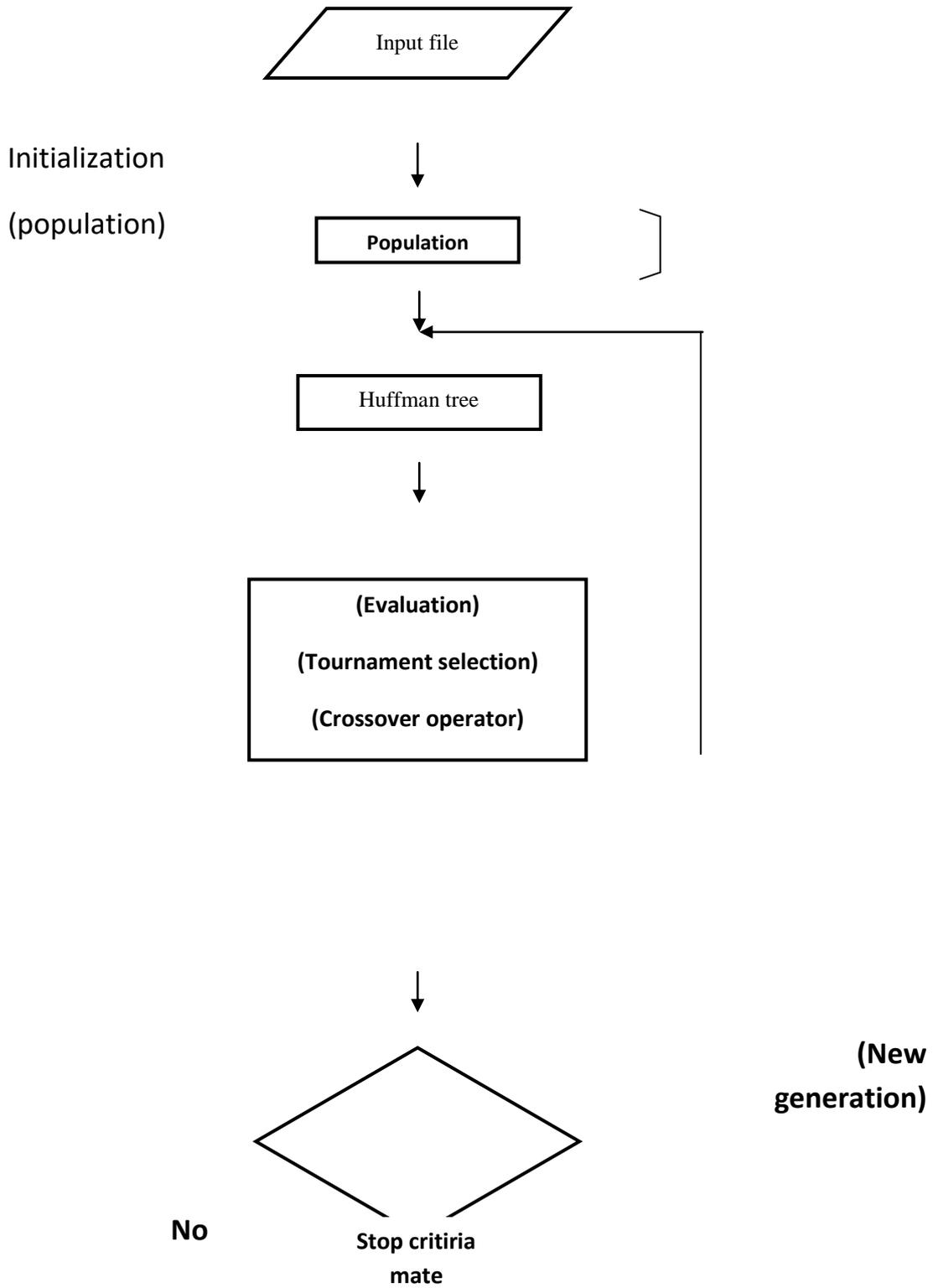


Figure (4-1) Block Daigram of Proposed System

The following flowchart in Figure (4-2) illustrates the direction of the system's work:



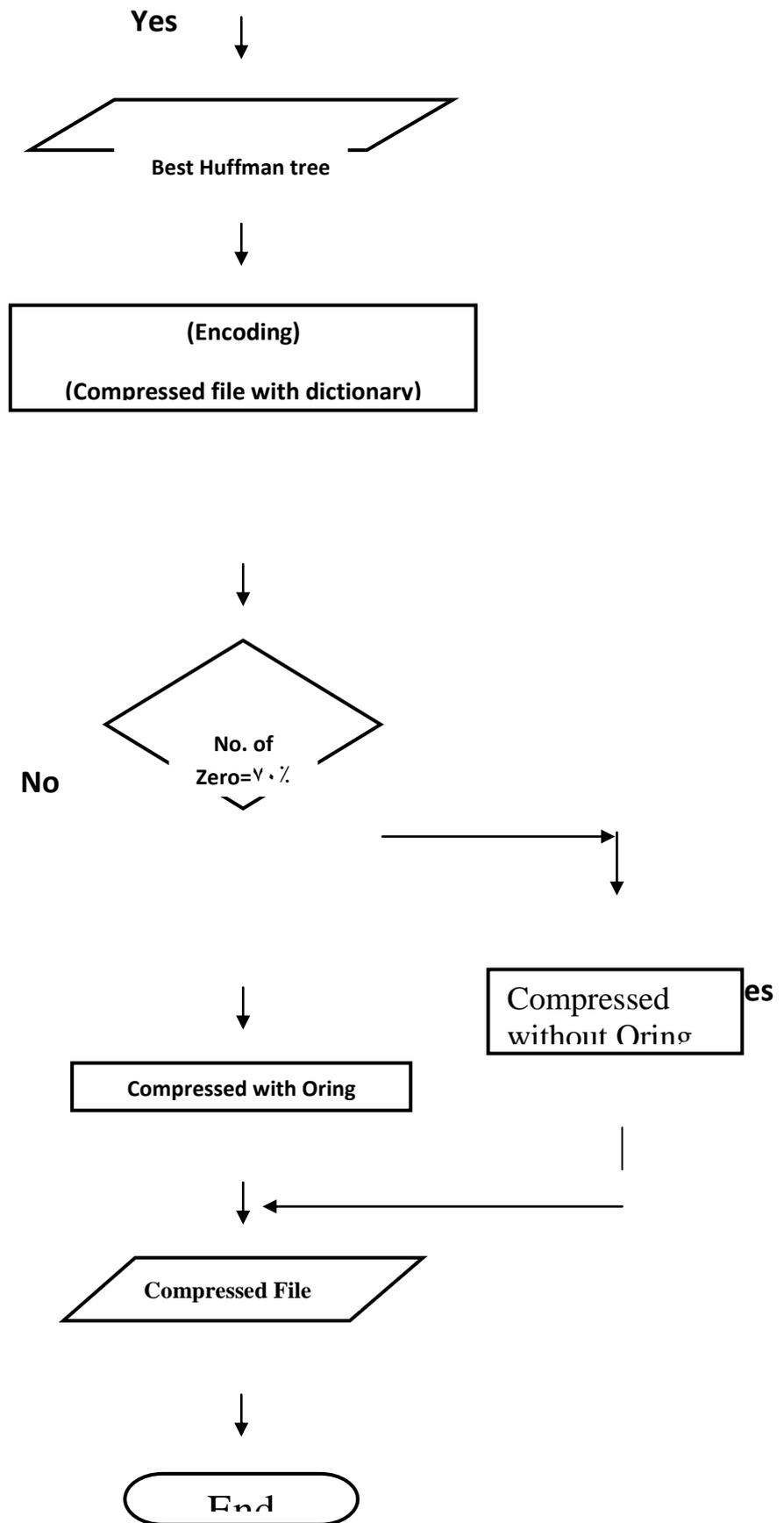


Figure (4-2) The General Flowchart for the Proposed System

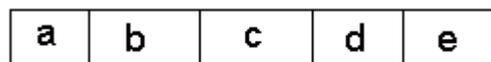
4.2 *Creating the Individual & the Initial Population*

First, the text is wanted to be compressed and find the best Huffman tree for it is inputted from the notepad as a file. It has used text files which are read from the proposed system depending on the text; the initial population is created, where an individual is produced and from this individual (chromosome) the other individuals of the population are found randomly.

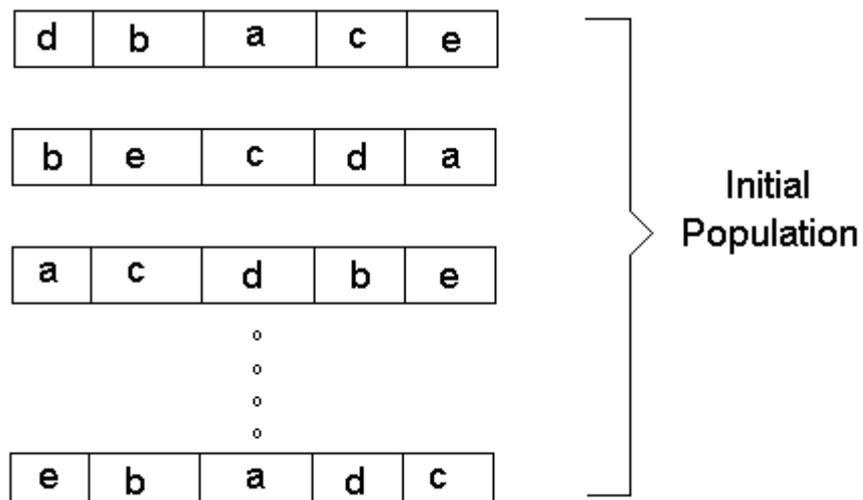
For example:

Text: abccddeeee

Produces the individual (chromosome)



:And from this individual the initial population is produced randomly as follows



The chromosome (individual) length depends on the variety in the text of the inputted characters. The frequency for each character is computed and then the probability for each character is found.

For example:

, c= 0.2 , d= 0.2 , e= 0.4 , b= 0.1 a= 0.1

The followed algorithm used to produce the individuals (chromosomes) and creating the initial population is:

- 1- Determining individual length.
- 2- Determining the population size.
- 3- Reading the input file and determining the first individual in the population.
- 4- Generating random function to exchange the Genes locations and probabilities in order to create the rest of the population.

4.3 Building Huffman Tree

After generating the initial population, the Huffman tree is built for each individual in the population depending on their probabilities. The codeword is determined from the tree for each gene, for example in Figure (4-3):

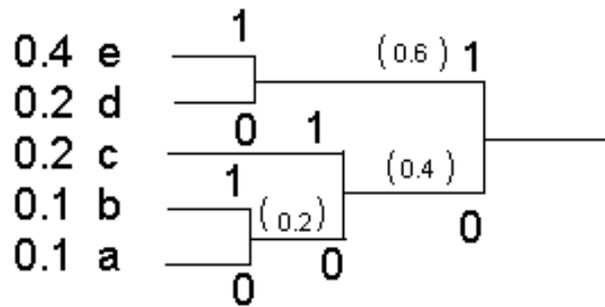


Figure (ε-ϣ) Huffman Codes

Where is:

e=11 d=10 , c=01 , b=001 , a=000 ,

The followed algorithm used for building the Huffman tree is:

- 1- Computing the probability for each character in the text.
- 2- Finding the two least probabilities according to the priority.
- 3- Giving value '0' for character with the biggest value from step 2 and value '1' for the other, this is on the condition that the two characters have two different values. In the condition of having the same value, the upper (partial branch) gets value '0' and the lower gets value '1'.
- 4- Sum the two probabilities.
- 5- The operation is repeated from step 2 to step 4 until the sum of two chosen probabilities is equal to '1'.
- 6- Retrieve codeword for each character from right to left (probability).

Where the Huffman tree represents the dictionary that will associate the compressed file and is used in the decoding operation.

4.4 Evaluation

Every individual is evaluated by computing the fitness function. The fitness function here is the variance that depends on computing the average size for each individual.

$$A = \sum_{i=1}^n (P_i * a_i) \quad , \quad \dots\dots\dots(\xi-1)$$

$$\text{Variance} = \sum_{i=1}^n p_i (a_i - A)^2 \quad , \quad \dots\dots\dots(\xi-2)$$

where,

p: probability

a_i: number of bits(length code word)

A: average size

For example:

The average size for the previous example shown in Figure (ξ-3) is computed as follows:

$$0.4 * 2 + 0.2 * 2 + 0.2 * 2 + 0.1 * 3 + 0.1 * 3 = 2.2$$

Then the variance is computed as follows:

$$0.4(\gamma - \gamma_0)^2 + 0.2(\gamma - \gamma_0)^2 + 0.2(\gamma - \gamma_0)^2 + 0.1(\gamma - \gamma_0)^2 + 0.1(\gamma - \gamma_0)^2 = 0.16$$

4.5 Tournament Selection

After the evaluation operation has been made, the two parents that will take place in the crossover operation are selected randomly.

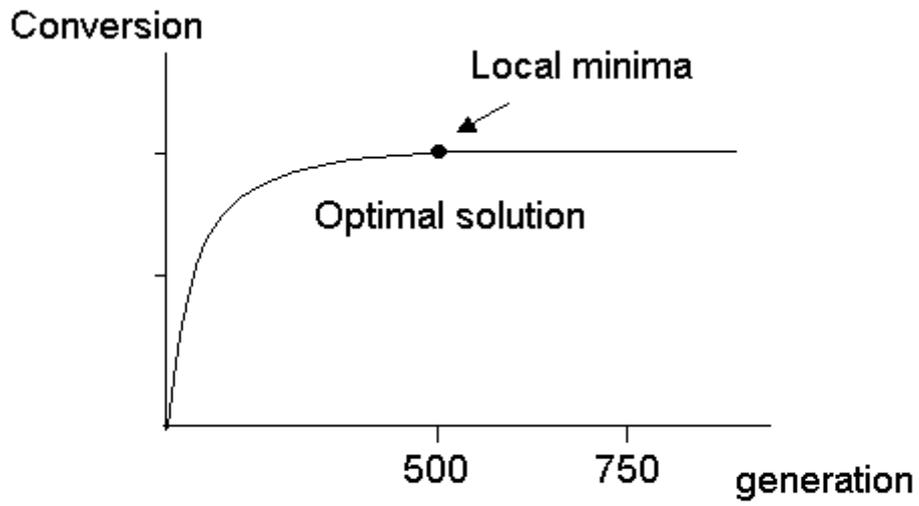
Two randomly individuals are chosen to produce the subpopulation; the individual with the least variance is selected to be the first parent. The operation is repeated to get the second parent according to the same parameters above.

The followed algorithm used in the proposed system is:

- 1- Finding the subpopulation by selecting two individuals randomly.
- 2- Selecting the individual with the best fitness and saving it; this is one of the two parents taking part in the crossover operation.
- 3- The operation from step 1 to step 2 is repeated to get the second selected parent.

Taking into consideration the subpopulation size that has an effective role in reaching the optimal solution, consider from Figure (4-4) that if the size of the subpopulation is more than 3, this will lead to premature conversion to a solution that may not be the optimal solution; this condition leads the system to local minima.

In the proposed system it has chosen the size of subpopulation equal to '3', in order to avoid the premature conversion.



Premature conversion

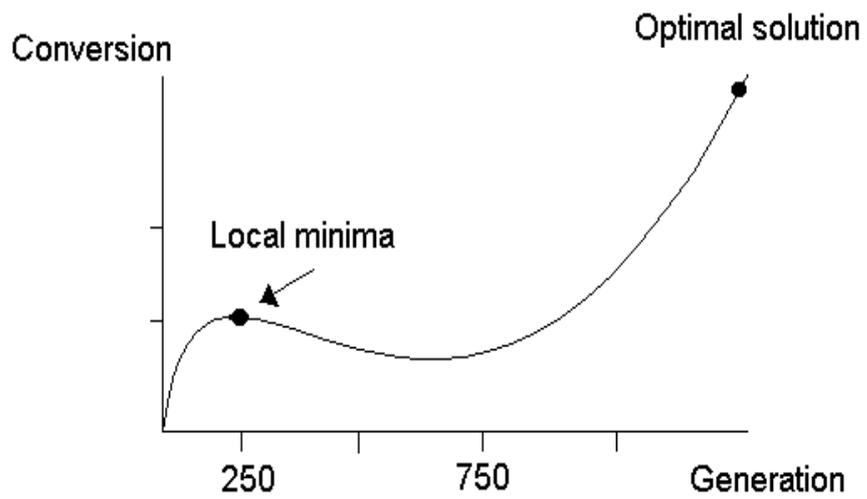


Figure (4-4) Conversion to the Optimal Solution

In general, the Figure (ξ-ο) illustrates the selection operation of individuals that share in each of the crossover, mutation, and the evaluation operation in order to obtain the offspring.

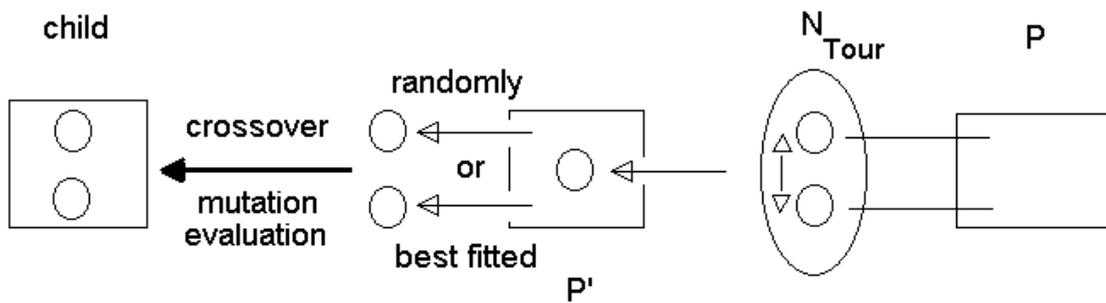


Figure (ξ-ο) Tournament Selection

Where is,

N_{tour} : subpopulation

P' : partial population

P : old population

The Figure (4-7) illustrates an example of a binary tournament selection. The population consists of a set of chromosomes whose genes are from the English alphabetic A.

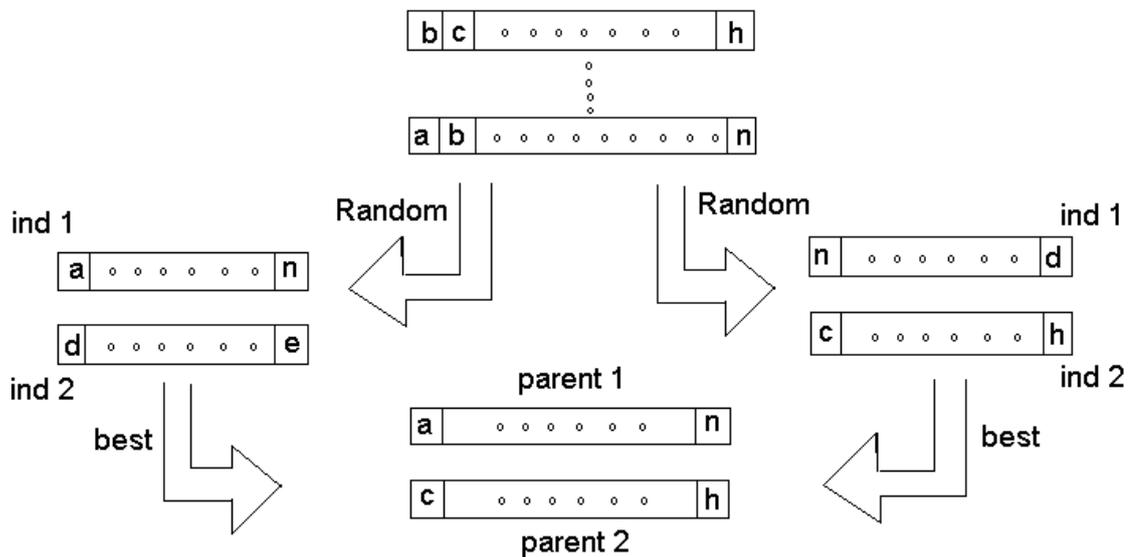


Figure (4-7) Binary Tournament Selection Example

4.6 Crossover Operation

In the proposed system, the cycle crossover (CX) has been chosen as one of the permutation crossover operators that mate the matching with the problem. This type of crossover gives a variety of individuals; in addition it avoids the conflict in the genes constructing the chromosome (individual), which is the most important property that must be available in the proposed system.

The general form for the cycle crossover is illustrated in Figure (4-8):

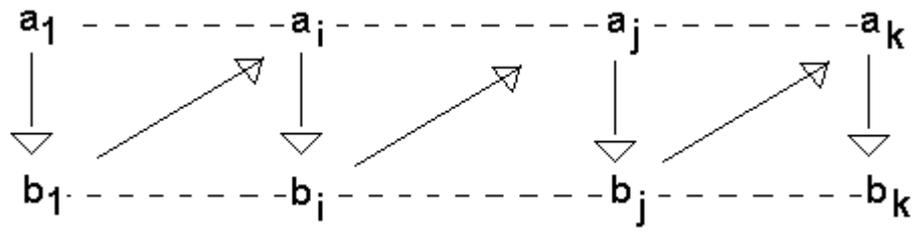


Figure (ξ-ν) Cycle Crossover

where $a_1, b_i, a_j, a_k, b_1, b_i, b_j, b_k$ represents the cycle elements(crossover cycle).
 The cycle is stopped when $a_1 = b_k$. The elements that have shared in the cycle keep the same location and the other elements exchange their locations.

The following example in Figure (ξ-λ) describes that:

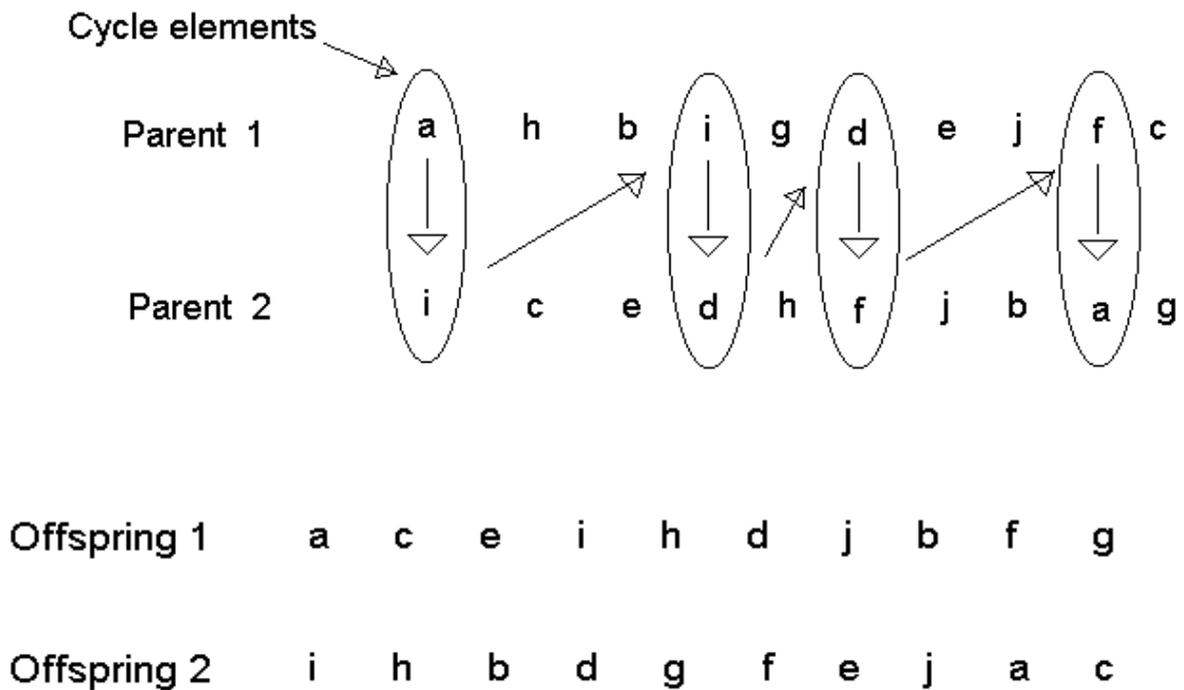


Figure (4-8) Crossover Operation

The followed algorithm of the cycle crossover is:

- 1- Reading the genes (characters) for the parents and their probabilities.
- 2- Selecting the first gene of both of the parents, and saving it in the variables 'one','two' respectively.
- 3- Comparing 'one' and 'two' variables, if 'one' and 'two' are equal print 'no crossover' and go to step 6. Otherwise, if 'one' and 'two' variables are not equal , get the similar gene from the first parent for the first gene of the second parent.
- 4- Selecting the similar copy gene of the second parent for the selected gene in step 3 (first parent), and saving it in 'two' variable.
- 5- Comparing 'one' and 'two' variables, if not equal go to step 3. Otherwise, if equal, exchange the genes that did not share in the cycle, according to the gene and the one corresponding to it.
- 6- End.

Mentioning that the crossover rate (P_c) is 1.

4.7 Mutation:

From computing the probabilities, it is determined whether there is mutation or not. The mutation probability (P_m) is 0.0001 , if the gene's probability is less than or equal to the P_m then mutation occurs at that gene. On condition mutation occurs, the gene's location that happened at it the mutation is exchanged with the succeeding gene. If the mutation occurs at the last gene, in this condition this gene's location is exchanged with the first gene.

Figure (4-9) illustrates the mutation operation:

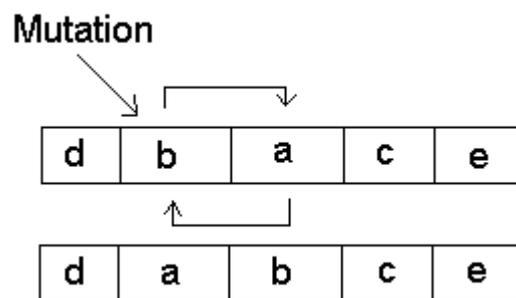


Figure (4-9) Mutation Operation

4.8 Evaluate Offspring and Replacement

The new offspring produced by the crossover operation after building the Huffman tree for each new individual are evaluated and the fitness function is computed.

Then the new offspring is compared with the worst individuals in the population (biggest variance), the offspring are exchanged with the worst individuals in case the offspring is better than or equal to the selected individuals in order to get the best variety in the population.

In other cases, no exchange happens. After finishing the evaluation operation, there will be a new selection and crossover operation again and continue in this way.

Noting that the proposed system had applied the Niching concept, where after beginning the search for the optimal solution, we notice that the system will find peaks of solutions moving between them to reach the optimal solution.

The benefit getting from applying the concept of Niching is providing a set of solution rather than only one solution, as shown in Figure (٤-١٠).

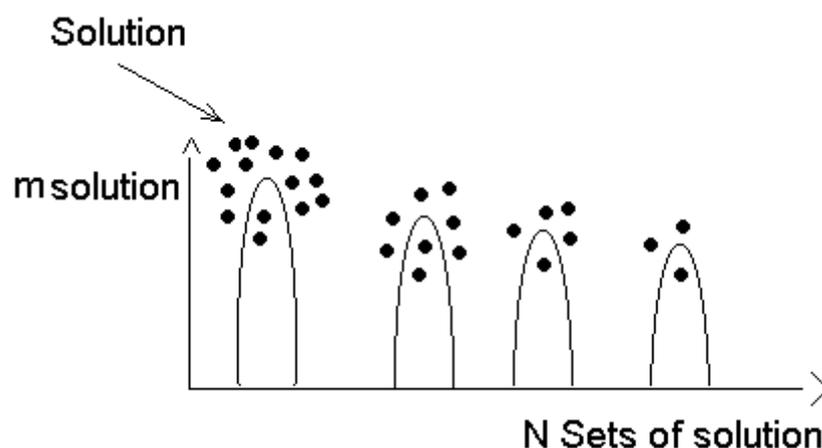


Figure (٤-١٠) Niching Concept (Subpopulation)

4.9 Stop Criteria

The operation including the crossover operation and generating a new population is continued for some generations (may be thousands of times).

After this, we will consider that one of the individuals with the best fitness (the least variance) prevalent in the population. The rate of prevailing may reach 90% or more and this rate is one of the stop criteria; the other stop criteria are the number of cycles (generations) that can reach 1000 generation.

As a result, the best individual is the most prevalent in the last population.

4.10 Compressed Operation

After finding the best Huffman tree (Dictionary), the text is coded according to the dictionary that has been obtained and according to the codeword for each character in the text. The result is a sequence of '0' and '1' saved in a file.

The number of zeros are tested, if the number increases 50% then the file is inputted to another additional compression algorithm working on binary numbers; this method is Oring bits. This method rises from the compression rate.

After the Oring bits method is applied on the resulting file, the binary numbers are divided to blocks, each block includes eight bits (byte) and is converted to an integer number and saved in byte. The file is sent as a compressed file to other side, mentioning that the dictionary will be associated with the text.

Example :

If we have the following text (file size is 100 byte):

aaaaaaaaaaaaaaaaaabbcccccccccccccccccccccccccccccccccccc

cc

After building the following Huffman tree, as shown in Figure (4-11):

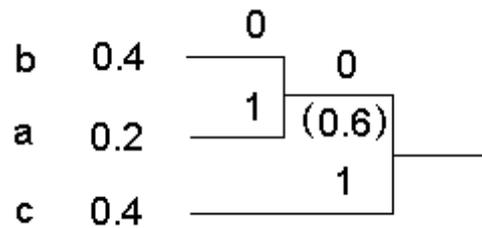


Figure (4-11) Huffman Tree

The codeword for each gene is :

c = 0 , a = 10 , b = 11 .

The resulted file is :-

```

11111111|11111111|11111111|11111111|11111111|.....|.....
|.....|.....|.....|.....|.....|.....|.....
|.....|11111111|11111111|11111111|11111111|

```

Noticing from the text that the number of zeros has a high rate, and by that the additional compress using Oring bits is useful, where the compression rate increases from 71% to 77% and as following:-

$L^1 = 01010101 | 01010101 | 01010101 | 01010101 | 01010101 | 11111111 | 11111111$
 $111 | 11111111 | 11111111 | 11111111$
 Oring $L^1 = 11111... |1 | 1111$
 $L^2 = 11111... |1 | 1111$
 Oring $L^2 = 111$

Then Oring L^2, L^2, L^1 is sent with the dictionary through the transmission channel and as follows (file of byte) :-

3 80 80 4 0 99 98 97 16.
 1 31 10 80 200 200 200 200 200 248 80 80

The first byte represents the number of genes that constructs the individual (chromosome length). The second byte and third byte represent the number of bits that represent text length where:

Text length = second byte * 200 + third byte, then the text length is equal to 16.

The next n bytes (n=number of genes) represents the genes (characters). The other next n bytes represent the codeword for each gene (dictionary). while the rest of bytes represent the coding of the text.

4.11 *Decompressed Operation*

The first thing done after receiving the compressed file is determining whether the Oring bits have been used or not.

The genes and the dictionary are distinguished from the compressed text, then the text is read bit by bit in a sequential order and is compared with the dictionary in order to get the original text.

In case the file is not inputted to the Oring bits method, then the integer values are converted to the corresponding binary values. So the file is now converted to a sequence of binary numbers (0 and 1) which represent the text.

On the other hand, if the Oring bits are used, then the rest of the bytes (text bytes) in the compressed file represent the Oring bits bytes. In this condition the decompression operation begins from the last byte which represents the Oring L^y .

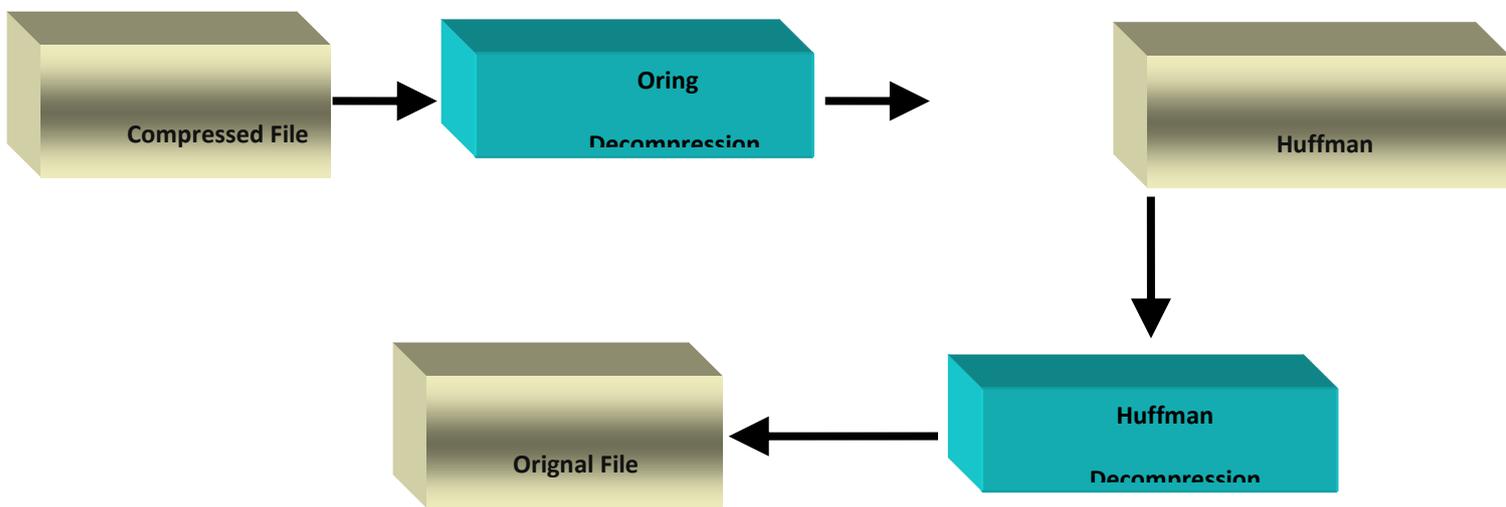
The Oring L^y is converted to binary values. From the number of one's in the binary value we determine the number of bytes for L^y which are also converted to binary in order to know the number of one's which give the number of bytes of L^x also converted to a binary value.

After determining each of Oring L^y , L^x , L^z , now the Oring bits procedure is applied. As a result, we have a file as a sequence of 0 and 1, which is the same result of the first case without Oring bits.

The resulting file is the text. The text is read bit by bit in a sequential order and is compared with the dictionary in order to get the original text.

The block daigram below illsturated this proses :

Decompression Process



Figure(4-11) Decompression process

o.1 Result Analysis

After building the Huffman tree for the individuals produced from

the genetic algorithm, we have to use a measurement for finding the best tree. In the beginning the Entropy was used to find the best tree. It was found out that all the trees have the same entropy because it depends on the probabilities of the genes that are the same in all the individuals and as a result cannot be used as a measurement for finding the best tree.

Then the variance was used for finding the best Huffman tree. It was found that there was a variety in the variance values that depend on the text itself; and the variety may be large or small. The benefit of the variance depends on the average size of the tree; we will get on codewords with approximate lengths and this will help in the transmission operation where the bandwidth is restricted in the transmission channel and finally making use of the buffer size and transmission time.

It has been noticed that several trees may have the same number of bits constructing them but it has different variance values and the best variance(best tree) is selected which have the codewords with approximate lengths.

These several examples illustrates the individuals and their codewords and also shows the role of the variance in choosing the best Huffman tree.

At first, the number of generation is determined as one of the stop criteria. Then the desired text is inputted as a file and the number of individuals of the initial population are also determined.

Example 1:

The number of generation is 100 and the number of individuals is 10.

abccddeeeeabccddeeeeabccddeeeeabccddeeee.....

The file size is 100 byte. The initial population is generated as:-

File size :100

Generation : 1

Ind : 1 Allcode = 1 4

Variance = 1.360

..... 1.100 [a]

..... 1.100 [b]

... 1.200 [c]

.. 1.200 [d]

. 1.400 [e]

abcde

ind : 2 Allcode = 1 4

variance = 1.360

..... 1.100 [a]

... 1.200 [d]

. 1.400 [e]

.. 1.200 [c]

..... 1.100 [b]

adecb

ind : 3 Allcode = 1 4

variance = 1.360

... 1.200 [c]

0010 0.1000 [a]

01 0.2000 [d]

0011 0.1000 [b]

1 0.4000 [e]

cadbe

The generation operation of individuals for the initial population continues and new individuals are generated with different number of bits for the generated tree that can be better or worse from the previous trees and as follows :

Ind: 7

allcode: 13

Variance: 0.960

000 0.2000 [d]

1 0.4000 [e]

001 0.2000 [c]

010 0.1000 [b]

011 0.1000 [a]

decba

After making the crossover operation between the individuals for some generations, a new individual may get produced carrying better characteristics

from the parents and with this it will be the individual with the least variance.

The following example illustrates that:

allcode = 12

Variance = 0.160

00 0.400 [e]

100 0.100 [b]

11 0.200 [d]

01 0.200 [c]

101 0.100 [a]

ebdca

Another individual may be produced carrying the characteristics of the previous individual and this individual may be prevalent in the population after some generations and by that it achieves the best individual in the population.

This example illustrates that :-

Best tree : ^

Variance : 0.160

00 0.400 [e]

100 0.100 [a]

۱۱ ۰.۲۰۰ [c]

۰۱ ۰.۲۰۰ [d]

۱۰۱ ۰.۱۰۰ [b]

eacdb

Table (۰-۱) shows the possible variance and Allcode (number of bits for the dictionary) for all the population's individuals for the previous text:-

Table (۰-۱) Allcode of Example ۱.

Individual	Allcode	Variance
abcde	۱۴	۱.۳۶۰
decba	۱۳	۰.۹۶۰
ebdca	۱۲	۰.۱۶۰

The relation between the variance and Allcode (number of bits for the dictionary) is shown in Figure(۰-۱):

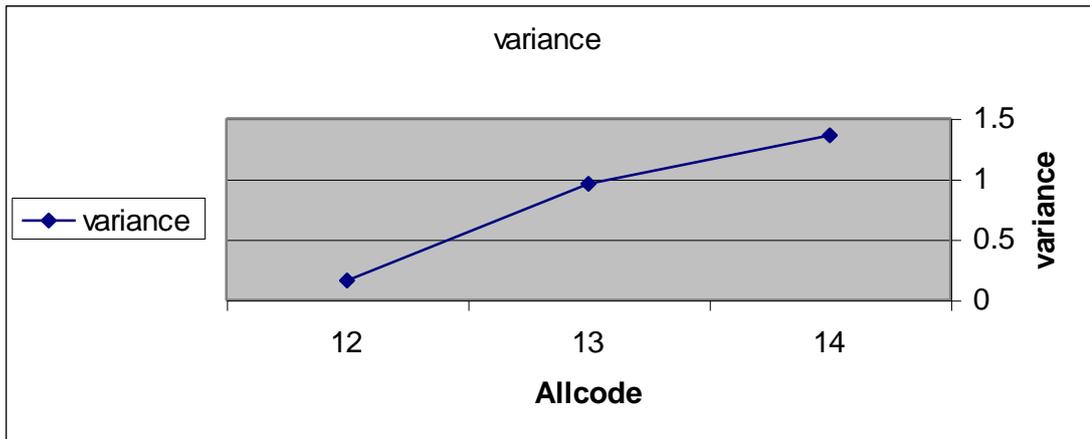


Figure (9-1) Variance of Example 1.

where the allcode represents the number of bits that the Huffman tree consists of. Noting that the number of the individuals of the initial population has a large and obvious effect on the results where the increase in the number of individuals of the initial population carryout increase in the crossover operation and which finally leads to a variety in the individuals in order to get new individuals and by that avoiding not obtaining a crossover operation that may happen from the fast prevalence for a particular individual.

In addition to the increase in the number of individuals of the initial population, the increase in the variety of the genes of the chromosome (increasing individual's length) leads to the variety in the individuals that will lead to the increase in the crossover operation and finally can obtain the best individual.

Example 2:

abccddeeeeqwrrttyyyynmlloo.....

Where the file size is 338 byte.

The variance table (0-2) that is obtained after 300 generation is :-

Table (0-2) Allcode of Example 2.

Individual	Allcode	Variance
teydorclmnqwba	04	0.213
tyeqorclabnwmd	00	0.367
abcdeqwrtynmlo	06	0.521
endmylabwcrotq	07	0.675
mcrloaytdwqnbe	07	0.828
elnrtwbmaoqcdy	08	0.982

As it is noticed from the Table (0-2), the same allcode value gives two different variance values. In this case, the structure of the tree has an effect on the codeword length and finally has an effect on the variance.

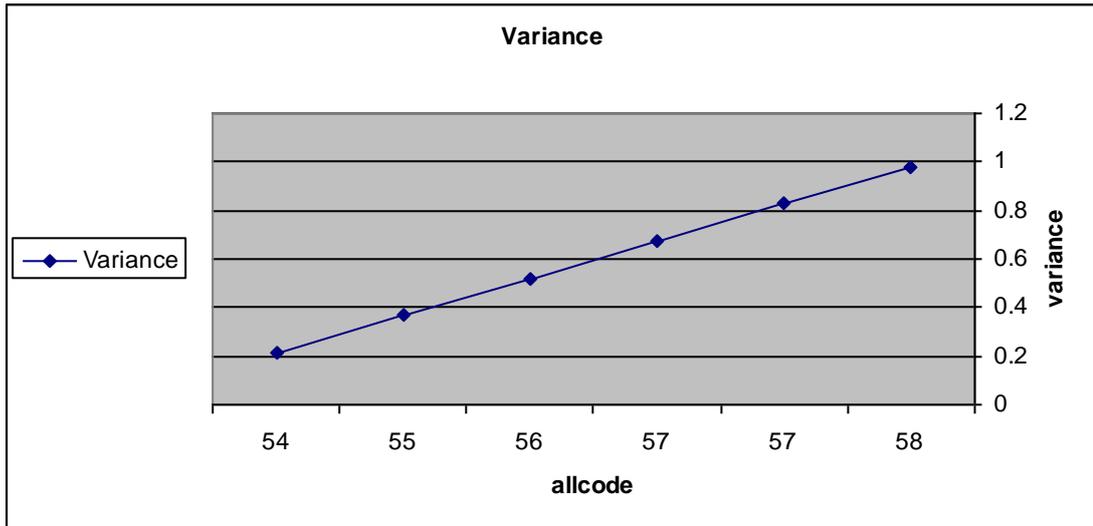


Figure (9-2) Variance of Example 2.

The randomize plays a large role in obtaining the results and the speed
in reaching the optimal result.

Example 3:

The boy is very clever and happy

Where the variety in the variance that can be obtained is shown in table (9-3):

Table (9-3) Allcode of Example 3.

Individual	Allcode	Variance
ahvdecynrtpoilsb	66	0.410

cedsrbpinvhyalto	٦٧	٠.٥٦٤
theboyisvrclandp	٦٨	٠.٧١٧
yvplronchtasdieb	٦٩	١.٠٢٥

Noting that after ٣٠٠ generations we can obtain the best variance that is

٠.٤١

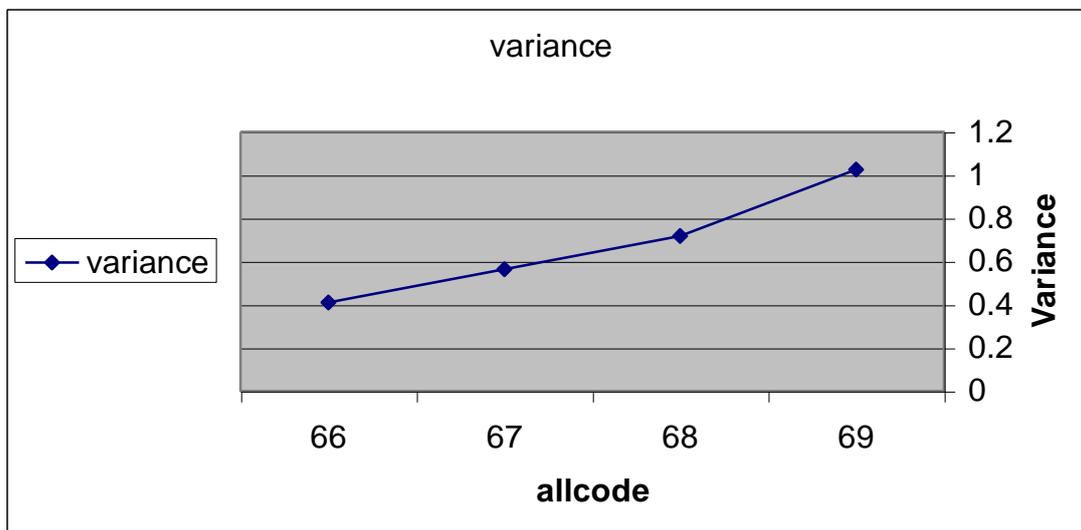


Figure (٤-٣) Variance of Example ٣.

The proposed system applied on a set of files has been chosen because the variety in the variance value in order to clear all the possible situations. It has been showed the effect of each of the file size, chromosome length, and probability of genes on the compression ratio as showed in the following tables. The Oring bits have been applied on a set of these files showing its effect in increasing the compression ratio.

Table (٥-٤) illustrates a set of examples taking into consideration different file sizes, while the chromosome length, population size, and generation number are fixed.

Table (٥-٤) Effect of Different File Size

	File name	File size	Chrom. Length	Population size	Generation No.	Compression Ratio
١	File ١	١٥٠B	٥	١٥	٢٠٠	٦٣.٣٣٣%
٢	File ٢	٣٠٠B	٥	١٥	٢٠٠	٦٨.٠٠٠%
٣	File ٣	٦٠٠B	٥	١٥	٢٠٠	٧٠.٣٣٣%

Table (٥-٥) illustrates a set of examples taking into consideration different chromosome lengths, while the other parameters are fixed.

Table (٥-٥) Effect of Different Chromosome Length

	File name	File size	Chrom. Length	Population size	Generation No.	Compression Ratio
١	File ٣	٦٠٠B	٥	٢٥	٣٠٠	٧٠.٣٣٣%
٢	File ٦	٦٠٠B	١٠	٢٥	٣٠٠	٥٦.٥٠٠%
٣	File ١١	٦٠٠B	٢٠	٢٥	٣٠٠	٤٠.٣٣٣%

Table (٥-٦) illustrates the effect of the probability of genes on the compression ratio. The increment in the probability of genes gives the gene the shortest codeword and as a result increases the compression ratio.

Table (٥- ٦) Effect of Probability of Gene

	File name	File size	Chrom. Length	Population size	Generation No.	Compression Ratio
١	File ١	١٥٠B	٥	١٥	٢٠٠	٦٣.٣٣٣%
٢	File ١٢	١٠٠B	٥	١٥	٢٠٠	٧١.٠٠٠%
٣	File ١٠	٩٠B	٥	١٥	٢٠٠	٧١.١١١%

After applying the Oring bit algorithm on some of the above files, taking into consideration that the ratio of zeros in the file must be above ٦٥% in order to get on good results.

The files that satisfied the condition and gave good results are illustrated in Table (٥-٧).

Table (٥- ٧) Results of Oring Files

	File name	File size	Chrom. Length	Compression ratio without Oring	Compression ratio with Oring
١	File ٦	٦٠٠B	١٠	٥٦.٥٠٠ %	٦٢.١٦٧ %
٢	File ١٠	٩٠B	٥	٧٠.١٠٣ %	٧٧.٣٢٠ %
٣	File ٩	٨٨٤B	٣	٧٩.١٨٦ %	٨٤.٦١٥ %

Also the proposed system has been applied on segments from DNA file taking different file sizes and results are viewed as follows:

Table (٥- ٨) Results of DNA Files

File name	File size	Compression ratio
DNA ١	٧٢ B	٥٩.٧٢٢٪
DNA ٢	١٥٣ B	٦٧.٣٢٠٪
DNA ٣	٣٠٠ B	٧١.٣٣٣٪
DNA ٤	٥٤٢ B	٧٢.٨٧٨٪

It has been noticed that there was not a large variety in the individuals because of the small number of genes (only four genes). It has also been noticed that it was difficult to get good results from applying the Oring on the DNA file because there was no occurrence of a large sequence of single genes.

٥.٢ Conclusions

١. The subpopulation size has an effective role in reaching to the optimal solution; if the size of the subpopulation is more than ٣, this will lead to premature conversion to a solution that may not be the optimal solution.
٢. The same allcode value gives two different variance values. In this case, the structure of the tree has an effect on the codeword length and finally has an effect on the variance.
٣. The number of individuals (population size) has an effect on obtaining the best tree where the small population size leads to a small variety and as a result leads to small crossover operations between the individuals because of getting the same first gene. As a result, the best individual (tree) will not be obtained.
٤. The number of generations has an effect on reaching to the best Huffman tree where continuing in selecting various individuals gives a larger chance in selecting the best possible individual.

- . In order to get on the best compression for the Oring method when merging the Huffman compression method and Oring, the character with the largest probability takes the value ‘0’ instead of the value ‘1’ in building the Huffman tree.
- ⌘. The chromosome length has an effect on the compression ratio. Whenever the chromosome length was longer, the compression ratio is smaller. This happens because the increase in the size of the Huffman tree causes as a result an increase in the number of bytes transmitted through the transmission channel.

◦.Ⓝ Future Works

Ⓛ- Applying the meta -Genetic Algorithm on Genetic Algorithm

(Optimization of GA).

A meta-genetic algorithm has been used to optimize the genetic algorithm for cell placement. The three parameters optimized are the crossover rate, inversion rate and mutation rate. The meta-genetic algorithm is itself a genetic optimization process which runs the genetic algorithm to solve a placement problem and manipulates the genetic parameters to optimize the fitness of the genetic algorithm. The individuals in the population of meta-genetic algorithm consist of three integers in the range $[0, 255]$, representing the mutation rate, inversion rate, and crossover rate for the genetic algorithm.

Ⓜ- Using Breeder Genetic Algorithm (BGA):-

BGA represents a class of random optimization techniques gleaned from the science of population genetics, which have proved their ability to solve hard optimization problems with continuous parameters. BGA which can be seen as a recombination between Evaluation strategies (ES) and Genetic Algorithm (GA), uses truncation selection which is very similar to the (μ, λ) strategy in ESs and the search process is mainly driven by recombination making BGAs very similar to GAs. It has been proven that BGAs can solve problems more efficiently than GAs due to the theoretical faster convergence to the optimum and they can, like GAs, be easily written in a parallel form.

۳- Applying the compression algorithm, the (Delta algorithm) on the research results for decreasing the cost of communication channels (reducing band width).

After finishing the compression operation and converting the sequence of binary numbers to integer numbers and before saving it as bytes, the delta algorithm is applying in order to reduce the value of the integer numbers and finally reducing the amplitude for the transmitted signal and this gives us a small band-width.

۴- The proposed system can be applied on images by dealing with an image as segments where Huffman algorithm is applied on each segment.

REFERENCES

- A. Joglekar & others, "**Genetic Algorithms and Their Use in the Design of Evolvable Hardware**", Fr. Conceicao Rodrigues College of Engineering, Bandra, Mumbai - ४०० ०००, IEEE, १९९८. [१]
- A. Khelifa , "**A Genetic Clustering for Image Segmentation** " , [२]
M.SC Thesis, April २००२.
- A. Armstrong and others, "**Genetic Algorithms for Image Compression** ", २००१ . [३]
- A. W. berger and others," **A Hybrid Coding Strategy for Optimized Test Data Compression**", University of Innsbruck, Austria, Proceedings IEEE International Test Conference, Charlotte, NC, USA, September ३० – October २, २००३. [४]
- [०] B. J. Schachter and others," **Some Experiments in Image Segmentation by Clustering of Local Feature Values** ", Pattern Recognition, ११, १, १९, १९८८.
- [१] B. D. Davison and others, "**Effect of Global Parallelsim on A Steady State GA** ", Proceedings of The Evolutionary Computing and Parallel Processing Workshop (GECCO'९९), Orlando, १९९९.
- D. Whitley, "**A Genetic Algorithm Tutorial** " , Computer [८]
Science Department Colorado State University, Fort Collins CO.,

- [۸] D. Salomon , “ **Data Compression the complete reference** “
 spring verlag Newyourk ,USA , ۱۹۹۸ .
- F. Hansson “ ,**Lempel-Ziv-Welch and Huffman compression** “ , [۹]
 ۲۰ January ۲۰۰۴ .
- G.B. Parker and others, “**Evolving Hexapod Gaits Using a Cyclic Genetic** [۱۰]
Algorithm ” , Department of Computer Science Indiana University
 Bloomington, ۱۹۹۸.
- Goldberg D. E. , “ **Genetic Algorithms in Search, Optimization** [۱۱]
and Machine Learning “ ,Addison-Wesley , ۱۹۸۹ .
- G. Kempe, “**Computer Science Honours Research Report** [۱۲]
Compression and Computational Gene Finding”, ۱ November
 ۲۰۰۲.
- G. E. Blelloch “**Introduction to Data Compression**” Computer [۱۳]
 Science Department Carnegie Mellon University , October ۱۶,
 ۲۰۰۱.
- H. J. Martinez and others, ” **Evolution of Cellular Automata** [۱۴]
for Digital Image Compression”, de Ingenieria Universidad
 Central de Venezuela Caracas, Venezuela, ۲۰۰۶.
- [۱۵] From Wikipedia, the free encyclopedia “**Huffman coding** “,
 GNU Free Documentation license, January ۱۹, ۱۹۹۶.

I. Harvey ,” **The Microbial Genetic Algorithm** “,School of [16]
Cognitive and Computing Sciences ,University of Sussex ,
Brighton BN1 9QH, UK, inmanh@cogs.susx.ac.uk., 1996.

J. Socha, “**Efficient Meta-Heuristics for Lossless Data** [17]
Compression”, School of Computer Science University of
Waterloo, Waterloo, Ontario, 2003.

J. M. Pullen,” **Data Compression, Security Principles** [18]
Integrity, Appropriate Use “, 2/3/03 © 2003.

[19] J. Smith & others, “ **Replacement Strategies in Steady State**
Genetic Algorithms : Static Environment “, Foundation of
Genetic Algorithms V, 219, 1998.

[20] M. Obitko ,“**Introduction to Genetic Algorithms with Java**
Applets”, 1998.

M. Sasikumar , “**Genetic Algorithms** “,NCST, Bombay. [21]

P. M. Elizabeth and others, “**Genetic Algorithm for VLSI** [22]
Design ,layout & test Automation “, 1999, prentice Hall ,
PTR USA.

[23] P. Buseti ,“**Genetic algorithms overview**”, 1990.

R. C. Gonzalez and others, “**Digital Image Processing**“, [24]
University of Tennessee, 1992.

R. He ,”**Indexing Compressed Text**”, A thesis, Waterloo, [25]
Ontario, Canada, 2003

R. Müller , “**Image Compression**” Part II: Image Processing [26]
Computer Graphics and Image Processing , Winter Semester
2003/04.

R. A. Watson and others, “**Recombination Without** [27]
Respect: Schema Combination and Disruption in Genetic
Algorithm Crossover “, *Volen Center for Complex Systems,*
Brandeis University, Waltham, 2000.

Sheila Horan, “**Data Compression Statistics and Implications**”, [28]
New Mexico State University, chapter 27- Data compression, 1997.

S. W. Smith, A sample chapter from: **The Scientist and** [29]
Engineer's Guide to Digital Signal Processing”, 1997 .

[30] T. A. Abbas Al- Asadi “**A Hybrid Algorithm For**
Images Compression” , Ph.D. Thesis, Computer Science ,
University of Technology, 2004 .

[31] T. H. Al-Kafaji ”**Improving Document Retrieval Using Genetic**
Algorithm” , M.SC. Thesis, September, 2000 .

U. Aickelin , “**Solving Multiple-Choice Problems with Genetic** [32]

Algorithms”, School of Computer Science University of

Nottingham UK uxa@cs.nott.ac.uk.

W. M. Spears “**Adapting Crossover in a Genetic Algorithm** “ , [33]

Washington, D.C. USA , 1990.

W. KEONG NG,” **Lossless and Lossy Data Compression** “ [34]

,Nanyang Technological University ,Singapore, 1996.

W. J. Mason and others , “ **Optimal Earth Orbiting Satellite** [35]

Constellations Via a Pareto Genetic Algorithm“ Department of

Aeronautical and Astronautical Engineering University of Illinois at

Urbana-Champaign Urbana, 1998.

From Wikipedia, the free encyclopedia,” **Genetic algorithm**” , [36]

GNU Free Documentation license, June 30, 2000.

REFERENCES