

**Implementation of one or more data compression algorithms,  
analysis of the algorithms and presentation of the results**

Project report submitted in partial fulfillment of the requirement for the  
degree of Bachelor of Technology

in

**Computer Science and Engineering/Information Technology**

By

MEGHA GOYAL (131233)

Under the supervision of

**PROF. DR. S. P. GHRERA**

to



Department of Computer Science & Engineering and Information  
Technology

**Jaypee University of Information Technology Waknaghat, Solan-  
173234**

## Candidate's Declaration

I hereby declare that the work presented in this report entitled “Implementation of one or more data compression algorithms, experimental comparison and analysis of the algorithm(s) and presentation of the results” in partial fulfillment of the requirements for the award of the degree of **Bachelor of Technology in Computer Science and Engineering/Information Technology** submitted in the department of Computer Science & Engineering and Information Technology, Jaypee University of Information Technology Waknaghat is an authentic record of my own work carried out over a period from August 2016 to May 2017 under the supervision of **Dr. Satya Prakash Ghrera (Professor, Brig (Retd.) and Head, Dept. of CSE and IT)**.

The matter embodied in the report has not been submitted for the award of any other degree or diploma.

Megha Goyal

131233

This is to certify that the above statement made by the candidate is true to the best of my knowledge.

Dr. Satya Prakash Ghrera

Professor, Brig (Retd.) and Head, Dept. of CSE and IT

Dated:

## Acknowledgement

It is our privilege to express our sincerest regards to our project supervisor (**Prof. Dr. Satya Prakash Ghjera**), for their valuable inputs, able guidance, encouragement, whole-hearted cooperation and direction throughout the duration of our project.

We deeply express our sincere thanks to our Head of Department **Prof. Dr. Satya Prakash Ghjera** for encouraging and allowing us to present the project on the topic “Implementation of one or more data compression algorithms, analysis of the algorithms and presentation of the results “at our department premises for the partial fulfillment of the requirements leading to the award of B-Tech degree.

## Table of Contents

<b>1. INTRODUCTION .....</b>	<b>1</b>
1.1 Introduction .....	1
1.2 Problem Statement .....	4
1.3 Objective .....	5
1.4 Methodology .....	5
1.4.1 Arithmetic Coding .....	5
1.4.2 LZW .....	7
1.4.3 Run-Length Encoding .....	8
1.4.4 Deflate .....	9
1.4.5 Prediction by Partial Matching .....	9
<b>2. LITERATURE SURVEY.....</b>	<b>10</b>
2.1 Arithmetic Coding .....	10
2.1.1 Practical Implementation.....	13
2.1.2 Underflow.....	14
2.2 LZW .....	16
2.2 Run-Length Encoding .....	19
2.2 Deflate .....	21
2.2 Prediction by Partial Matching .....	22
<b>3. SYSTEM DEVELOPMENT .....</b>	<b>24</b>
3.1 Arithmetic Coding.....	24
3.1.1 Algorithm .....	24
3.1.2 Model Development .....	25
3.1.3 Analysis of arithmetic coding .....	27
3.2 LZW .....	29
3.2.1 Algorithm .....	29
3.2.2 Model Development .....	29
3.2.3 Analysis of LZW.....	31

3.3 Run-Length Encoding .....	34
3.3.1 Algorithm .....	34
3.3.2 Model Development .....	35
3.3.3 Analysis of run-length encoding .....	36
3.4 Deflate .....	37
3.4.1 Algorithm .....	37
3.4.2 Model Development .....	39
3.4.3 Analysis of deflate .....	39
3.5 Prediction by Partial Matching .....	41
3.5.1 Algorithm .....	41
3.5.2 Model Development .....	42
3.5.3 Analysis .....	42
<b>4. PERFORMANCE ANALYSIS</b> .....	<b>43</b>
4.1 Compression Ratio .....	43
4.1 Compression Speed .....	44
4.3 Compression Time .....	45
4.4 Results .....	45
<b>5. CONCLUSION</b> .....	<b>51</b>
5.1 Conclusion .....	51
5.2 Future Scope .....	52
5.3 Application Contribution .....	53
REFERENCES .....	55

## List of Figures

Fig 1-1: Data Compression techniques	1
Fig 1-2: System with typical processes for data compression. Arithmetic coding is normally the final stage, and the other stages can be modeled as a single data source $\Omega$ .	6
Fig 2-1: Arithmetic Data compression Flow Graph	12
Fig 2-2 (a) Encoding of "BILL GATES" (b) Decoding of "BILL GATES"	12
Fig 2-3: How the file will be processed into streams.	22
Fig 3-1: Flow chart of LZW compression Algorithm	32
Fig 3-2: Flow chart of LZW decompression Algorithm	33
Fig 3-3: Basic flow chart of Run-length Encoding Algorithm	35
Fig 3-4: Flow chart of Deflate Algorithm	40

## List of Graphs

Graph 1: Comparison of compressed file size of Sources.50MB	48
Graph 2: Comparison of compressed file size of English.50MB	48
Graph 3: Comparison of compressed file size of Pitches.50MB	49
Graph 4: Comparison of compressed file size of Proteins.50MB	49
Graph 5: Comparison of compressed file size of Dna.50MB	50
Graph 6: Comparison of compressed file size of Xml.50MB	50

## List of Tables

Table 1: Comparison between arithmetic and Huffman coding methodologies	16
Table 2: Results of Arithmetic Coding	45
Table 3: Results of LZW algorithm	46
Table 4: Results of Run-Length Encoding	46
Table 5: Results of Deflate algorithm	47
Table 6: Results of Prediction by partial matching algorithm (PPM)	47



## Abstract

Data Compression is a strategy for encoding decides that permits considerable diminishment in the aggregate number of bits to store or transmit a document. Transmission of large quantity of data cost more money. Hence choosing the best data compression algorithm is really important. In addition to different compression technologies and methodologies, selection of a good data compression tool is most important. There is a complete range of different data compression techniques available both online and offline working such that it becomes really difficult to choose which technique serves the best. In this report I represent five algorithms (Arithmetic Coding, LZW, Run-Length Encoding, Deflate and Prediction by Partial Matching) to compress and decompress the text data.

# Chapter-1

## INTRODUCTION

### 1.1 Introduction

Data compression is a procedure that lessens the information measure, expelling the extreme data and repetition. Why shorter information grouping is more reasonable? –the answer is straightforward it lessens the cost. Information compression is a typical necessity for the greater part of the electronic application [1]. Information compression has imperative application in the range of document stockpiling and circulated framework. Information compression is utilized as a part of sight and sound field, content records and database table. Information compression strategies can be grouped in a few ways. A standout amongst the most essential criteria of arrangement is whether the compression calculations expel some piece of information which can't be recouped amid decompression. The calculation which expels some piece of information is called lossy information compression. The lossy information compression calculation is typically utilize when a flawless consistency with the first information is redundant after decompression. Case of lossy information compression is compression of video or picture information. Lossless information compression is utilized as a part of content document, database tables and in therapeutic picture since law of directions. Different lossless information compression calculation have been proposed and utilized. Some of primary methods are Huffman Coding, Run Length Encoding, Arithmetic Encoding and Dictionary Based Encoding. In this report we analyze Arithmetic Encoding and Dictionary-based Algorithm and give examination between them as indicated by their exhibitions.

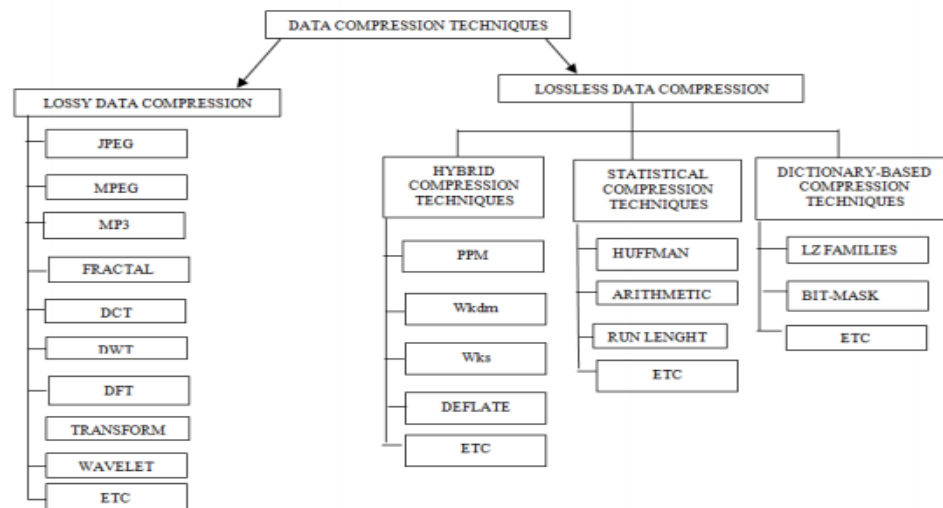


Fig 1-1: Data Compression techniques

Compression is utilized pretty much all over the place. Every one of the pictures you get on the web are compacted, commonly in the JPEG or GIF designs, most modems utilize compression, HDTV will be packed utilizing MPEG-2, and a few document frameworks naturally pack records when put away, and whatever remains of us do it by hand. The flawless thing about compression, as with alternate subjects we will cover in this course, is the calculations utilized as a part of this present reality make substantial utilization of a wide arrangement of algorithmic apparatuses, including sorting, hash tables, tries, and FFTs. Moreover, calculations with solid hypothetical establishments assume a basic part in true applications.

In this section we will utilize the bland term message for the articles we need to pack, which could be either documents or messages. The undertaking of compression comprises of two parts, an encoding calculation that takes a message and produces a "compacted" portrayal (ideally with less bits), and an unraveling calculation that recreates the first message or some estimation of it from the packed portrayal. These two parts are regularly complicatedly entwined since they both need to comprehend the common compacted portrayal. We recognize lossless calculations, which can remake the first message precisely from the packed message, and lossy calculations, which can just reproduce an estimate of the first message. Lossless calculations are normally utilized for content, and lossy for pictures and sound where a tiny bit of misfortune in determination is regularly imperceptible, or if nothing else satisfactory. Lossy is utilized as a part of a dynamic sense, be that as it may, and does not mean irregular lost pixels, but rather implies loss of an amount, for example, a recurrence segment, or maybe loss of clamor. For instance, one may surmise that lossy content compression would be inadmissible in light of the fact that they are envisioning missing or exchanged characters. Consider rather a framework that rephrased sentences into a more standard shape, or supplanted words with equivalent words so that the record can be better packed. In fact the compression would be lossy since the content has changed, however the "signifying" and clearness of the message may be completely kept up, or even moved forward. Actually Shrunk and White may contend that great written work is the specialty of lossy content compression.

Is there a lossless calculation that can pack all messages? There has been no less than one patent application that asserted to have the capacity to pack all records (messages)—. The patent application guaranteed that on the off chance that it was connected recursively, a record could be lessened to nothing. With a little thought you ought to persuade yourself this is impractical, in any event if the source messages can contain any piece grouping. We can see this by a basic checking contention. We should consider each of the 1000 piece messages, for instance. There are 21000 diverse messages we can send, every which should be unmistakably recognized by the decoder. It ought to be clear we can't speak to that a wide range of messages by sending 999 or less bits for every one of the messages — 999 bits would just permit us to send 2999 particular messages. In all actuality if any one message is abbreviated by a

calculation, then some other message should be stretched. You can confirm this practically speaking by running GZIP on a GIF file.

As mentioned in the introduction, coding is the job of taking probabilities for messages and generating bit strings based on these probabilities. How the probabilities are generated is part of the model component of the algorithm, which is discussed in Section 4.

In practice we typically use probabilities for parts of a larger message rather than for the complete message, e.g., each character or word in a text. To be consistent with the terminology in the previous section, we will consider each of these components a message on its own, and we will use the term message sequence for the larger message made up of these components. In general each little message can be of a different type and come from its own probability distribution. For example, when sending an image we might send a message specifying a color followed by messages specifying a frequency component for that color. Even the messages specifying the color might come from different probability distributions as the probability for particular colors might depend on the context.

Also compressed files are much more easily exchanged over the internet as they upload and download much faster. We require the capacity to reconstitute the first record from the compacted rendition whenever. Information compression is a system for encoding that decides that permits a considerable decrease in the aggregate number of bits to store or transmit a document.

## 1.2 Problem Statement

The principal issue of lossless compression is to disintegrate an informational index (for instance, a content record or a picture) into a succession of occasions, then to encode the occasions utilizing as few bits as would be prudent. The thought is to allocate short code words to more likely occasions and longer code words to less plausible occasions. Information can be compacted at whatever point a few occasions are more probable than others. Factual coding methods utilize appraisals of the probabilities of the occasions to allocate the code words. Given an arrangement of commonly particular occasions  $e_1, e_2, e_3, \dots, e_n$ , and a precise appraisal of the likelihood  $P$  of the occasions, Shannon [3] demonstrated that the smallest conceivable expected number of bits expected to encode an occasion is the entropy of  $P$ , signified by

$$H(P) = \sum_{i=1}^n -p\{e_i\} \log_2 p\{e_i\}$$

where  $p\{e_i\}$  is the likelihood that occasion  $e_i$  happens. An ideal code yields  $\log_2 p$  bits to encode an occasion whose likelihood of event is  $p$ . Immaculate number juggling codes provided with exact probabilities give ideal compression. In principle, number-crunching codes dole out one "code word" to every conceivable informational collection. The code words comprise of half-open subintervals of the half-open unit interval  $[0,1)$ , and are communicated

by determining enough bits to recognize the subinterval comparing to the genuine informational collection from all other conceivable subintervals. Shorter codes compare to bigger subintervals and hence more plausible information informational collections. By and by, the subinterval is refined incrementally utilizing the probabilities of the individual occasions, with bits being yield when they are known. Number juggling codes quite often give preferable compression over prefix codes, yet they do not have the immediate correspondence between the occasions in the information informational index and bits or gatherings of bits in the coded yield record.

A factual coder must work in conjunction with a modeler that gauges the likelihood of every conceivable occasion at each point in the coding. The likelihood show require not depict the procedure that produces the information; it just needs to give a likelihood circulation to the information things. The probabilities don't need to be especially exact, however the more precise they are, the better the compression will be. In the event that the probabilities are uncontrollably incorrect, the document may even be extended as opposed to compacted, yet the first information can at present be recouped. To acquire greatest compression of a record, we require both a decent likelihood display and a productive method for speaking to (or taking in) the likelihood show.

Lossless information compression is a methodology that allows the use of information compression estimations to pack the substance information promote more allows the exact extraordinary information to be revamped from the compacted information. This is in instead of the lossy information compression in which the cautious one of a kind information can't be reproduced from the compacted information. Since the majority of this present reality data has factual excess, thusly lossless data compression is conceivable. Case in point, In English content, the letter "a" is a great deal more basic than the letter 'z', and the likelihood that the letter "t" will be trailed by the letter "z" is little. So this sort of repetition can be evacuated utilizing lossless compression. Lossless compression techniques may be classified by kind of data they are intended to pack.

### 1.3 Objective

Our objective it to implement the arithmetic coding, LZW, run length encoding, deflate and prediction by partial matching and compare the results obtained to maximize the compression ratio and minimize the compression time.

### 1.4 Methodology

### 1.4.1 Arithmetic coding

Arithmetic coding is an information compression procedure that encodes information (the information string) by making a code string which speaks to a fragmentary interval on the number line in the vicinity of 0 and 1. The coding calculation is highly recursive; i.e., it works upon and encodes (deciphers) one information image for each emphasis or recursion. On every recursion, the calculation progressively parcels an interval of the number line in the vicinity of 0 and 1, and holds one of the allotments as the new interval. In this manner, the calculation progressively manages smaller intervals, and the code string, seen as a size, lies in each of the settled intervals. The information string is recouped by utilizing size correlations on the code string to reproduce how the encoder must have progressively apportioned and held each settled subinterval. Arithmetic coding varies significantly from the more natural compression coding systems, for example, prefix (Huffman) codes. Additionally, it ought not be mistaken for error control coding, whose purpose is to recognize and redress errors in PC operations.

There are many preferences for isolating the source modeling (probabilities estimation) furthermore, the coding forms [14, 25, 29, 38, 45, 51, 53]. For instance, it permits us to create complex compression plans without stressing over the points of interest in the coding calculation, and/or use them with different coding methods and implementations. The two processes can be separated in a complete system for arithmetic encoding and decoding. The coding part is responsible only for updating the intervals, i.e., the arithmetic encoder implements recursion, and the arithmetic decoder implements. The encoding/decoding processes use the probability distribution vectors as input, but do not change them in any manner. The source modeling part is responsible for choosing the distribution  $c_k$  that is used to encode/decode symbol  $s_k$ . It also shows that a delay of one data symbol before the source-modeling block guarantees that encoder and decoder use the same information to update  $c_k$ . Arithmetic coding simplifies considerably the implementation of systems because the vector  $c_k$  is used directly for coding. With Huffman coding, changes in probabilities require re-computing the optimal code, or using complex code updating techniques [9, 24, 26].

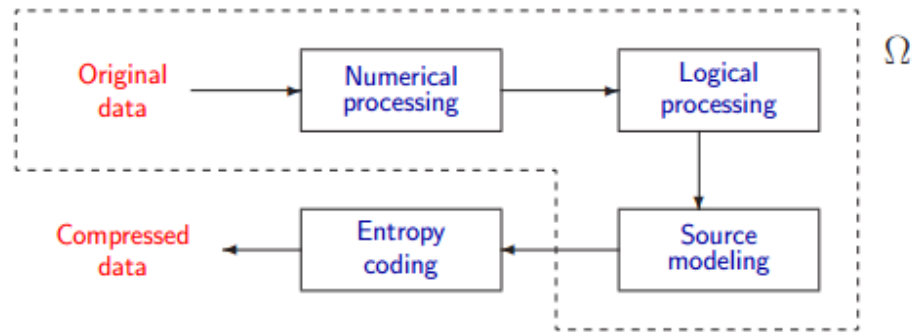


Fig 1-2: System with typical processes for data compression. Arithmetic coding is normally the final stage, and the other stages can be modeled as a single data source  $\Omega$ .

#### 1.4.2 LZW

LZW compression is the compression of a document into a littler record utilizing a table-based query calculation concocted by Abraham Lempel, Jacob Ziv, and Terry Welch. Two ordinarily utilized record designs in which LZV compression is utilized are the GIF picture arrange served from Web destinations and the TIFF picture organize. LZW compression is likewise appropriate for compacting content records.

A specific LZW compression calculation takes each info succession of bits of a given length (for instance, 12 bits) and makes a passage in a table (now and again called a "word reference" or "codebook") for that specific piece design, comprising of the example itself and a shorter code. As information is perused, any example that has been perused before results in the substitution of the shorter code, viably compacting the aggregate sum of contribution to something littler. Not at all like prior methodologies, known as LZ77 and LZ78, the LZW calculation includes the look-into table of codes as a feature of the packed document. The deciphering program that uncompresses the record can construct the table itself by utilizing the calculation as it procedures the encoded input.

As specified before, static coding plans require some information about the information before encoding happens. All inclusive coding plans, as LZW, don't require propel learning and can manufacture such information on-the-fly. LZW is the preminent system for universally useful information compression because of its effortlessness and flexibility. It is the premise of numerous PC utilities that claim to "twofold the limit of your hard drive". LZW compression utilizes a code table, with 4096 as a typical decision for the quantity of table passages. Codes 0-255 in the code table are constantly allotted to speak to single bytes from the information record. When encoding starts the code table contains just the initial 256 passages, with the rest

of the table being spaces. Compression is accomplished by utilizing codes 256 through 4095 to speak to successions of bytes. As the encoding proceeds with, LZW distinguishes reshaped groupings in the information, and adds them to the code table. Interpreting is accomplished by taking each code from the packed document, and deciphering it through the code table to discover what character or characters. The objective of word reference based displaying to infer an arrangement of expressions that can be utilized to financially speak to the message. Besides, since in a disconnected technique the expression table must be transmitted as a component of the packed message, the deduction plot utilized ought to permit a smaller encoding of the expression set. This last prerequisite does not have any significant bearing to incremental word reference based techniques.

Welch proposed an alteration of the more established LZ78 compression calculation. The proposed calculation, known as LZW, is a lexicon based compression calculation where the word reference is built as the info is prepared. The word reference  $D$  begins with all single-character strings. At every cycle, the following characters are filtered and the calculation coordinates the longest prefix of the info  $p$  that is in  $D$ , (i.e.,  $p \in D$  yet  $p \_ c \notin D$ , where  $c$  is the following character). The file of  $p$  in  $D$  is yielded and the string  $p \_ c$  is added to  $D$ . The calculation then reshapes the procedure on the rest of the information beginning with the following character  $c$ . This is appeared in calculation.

Decompression works similarly yet prepare files rather than characters: perusing in the following list, looking into the comparing section in  $D$ , and yielding the string  $s$ . It then takes the principal character  $c$  of  $s$  and affixes it to the past yield string  $s_0$ , embeddings  $s_0 \_ c$  into  $D$ . Decompressing is point by point in Algorithm 4 in Appendix B. Take note of that there is an extraordinary case in decompression: the information may contain the string  $c \_ p \_ c \_ p \_ c$ , where  $c \_ p$  is as of now in  $D$ . The compression calculation will coordinate  $c \_ p$  and embed  $c \_ p \_ c$  into  $D$ , then match  $c \_ p \_ c$  and yield its file. Be that as it may, at the beneficiary, the second list alludes to an unfilled section in  $D$ . Since we accept no blunders in the information, we realize that if the file focuses to the following accessible cell (where the new passage will be put), then we are in this exceptional case. Be that as it may, if the file does not indicate the following accessible cell then the information has been defiled and we have a translating blunder. The table utilized can be of settled size or it can develop progressively. The previous requires either some kind of deterministic expulsion approach or solidifying the lexicon when it turns out to be full. Developing the word reference progressively requires that the encoder and decoder develop the table in the meantime.

Regularly, the table is multiplied in size (and yields end up plainly one piece longer) when a section fills the table. This approach is most basic as it gives better compression, utilizing less bits from the get-go in the encoding procedure. Take note of that our figure is good with any lexicon administration conspire that is not subject to the request of the passages in the word reference (e.g., slightest as of late utilized).



The LZW calculation peruses the info characters from left to right while embeddings in  $D$  all substrings of the shape  $T[bm : bm+1]$ . Consequently the expressions of LZW are the substrings acquired by linking the pieces of  $T$  with the following character tailing them, together with every single conceivable substring of size one. The codeword of the expression  $T[bm : bm+1]$  is the whole number  $j_{j+m}$ , where  $j_j$  is the extent of the letter set  $_$ . In this way, the codewords of substrings don't change in LZW calculation. LZW utilizes covetous parsing too: the  $m$ th piece  $T_m$  is recursively characterized as the longest substring which is in  $D$  just before  $C$  peruses  $T[bm+1 ? 1]$ . Thus, no two expressions can be indistinguishable in the LZW calculation.

### 1.4.3 Run Length Encoding

Run length encoding can be found in various applications, for example, information exchange or picture putting away. It is a notable, simple and productive compression strategy in view of the supposition of long information successions without the change of substance. These arrangements can be depicted by their position and length of appearance. Executions utilizing devoted rationale are improved for parallel information handling. Here, pictures are moved in pieces of different pixels in parallel. A compression of these streams into a run length code requires an encoder with a parallel information. This run length encoder needs to pack the succession at least clock cycles to evade long hinder interims at the info. An equipment calculation playing out an elite run length encoding for paired pictures utilizing a parallel info.

### 1.4.4 Deflate

This determination characterizes a lossless compacted information arrange that packs information utilizing a mix of the LZ77 calculation and Huffman coding, with productivity equivalent to the best right now accessible universally useful compression techniques. The information can be delivered or expended, notwithstanding for a discretionarily long successively displayed input information stream, utilizing just a from the earlier limited measure of middle of the road stockpiling.

The motivation behind this particular is to characterize a lossless packed information organize that:

- Is autonomous of CPU sort, working framework, document framework, and character set, and subsequently can be utilized for trade;
- Can be created or devoured, notwithstanding for a self-assertively long successively exhibited input information stream, utilizing just a from the earlier limited measure of middle of the road stockpiling, and consequently can be utilized as a part of information correspondences or comparative structures, for example, Unix channels;

- Compresses information with productivity similar to the best at present accessible broadly useful compression techniques, specifically extensively superior to the "pack" program;
- Can be executed promptly in a way not secured by licenses, and consequently can be polished uninhibitedly;
- Is perfect with the document arrange delivered by the current generally utilized gzip utility, in that accommodating decompressors will have the capacity to peruse information created by the current gzip compressor.

#### 1.4.4 Prediction by Partial Matching

Prediction by Partial Matching (PPM) is a lossless compression calculation which reliably performs well on content compression benchmarks. The calculation is assessed on the Pizza corpus. PPM is worried with the principal errand of producing a likelihood conveyance for the expectation of the following character in an arrangement. Give us a chance to expect that we have a record, in which we can locate a since quite a while ago rehashed (at an irregular separation) string with the length of  $s$  bytes (characters). PPM needs to encode each character of this string independently (yet with high likelihood). PPM should dependably have a plausibility to encode the escape image (when another image in setting has showed up). For this situation, it is likely that the LZ77 calculation would pack better, as it encodes just the balance, the length, and the following character. Yet, in the event that we broaden the PPM letters in order with another image, which is equivalent to our string (with the length  $s$ ), we can encode this string significantly more effectively. Give us a chance to accept that a word reference is set of all since a long time ago rehashed strings. For this situation, our letters in order comprises of 256 characters (the normal PPM letters in order), the escape image, and images, which are equivalent to since quite a while ago rehashed strings from the word reference. We need to pre-ignore information to figure insights and find rehashed strings to develop the letters in order before we begin the PPM compression.

## Chapter-2

### LITERATURE SURVEY

#### 2.1 Arithmetic Coding

Data Compression systems are sorted by loss of information into two gatherings, in particular lossless information compression procedures and lossy information compression methods. Lossless calculations reproduce the first message precisely from the compacted message, and lossy calculations just remake a guess of the first message. Lossless calculations are regularly utilized for content and lossy for pictures and sound where a smidgen of misfortune in determination pitch or other is frequently imperceptible, or if nothing else acknowledged. Lossy is utilized as a part of a theoretical sense, notwithstanding, and does not mean arbitrary lost pixels, but rather implies loss of an amount, for example, recurrence segment, or maybe loss of commotion [4]. Lossy compression calculations are JPEG (Joint Photographic Expert Group), MPEG (Moving Picture Experts Group), MP (Media Player) and Fractal while Lossy coding methods incorporate DCT (Discrete Cosine Transform), DFT (Discrete Fourier Transform), DWT (Discrete Wavelet Transform). Lossless compression procedures are utilized to pack, of need, medicinal pictures, content and pictures saved for lawful reasons, PC executable records, among others. Lossy compression procedures recreate the first message with loss of some data. It is unrealistic to reproduce the first message utilizing the interpreting procedure, and is called irreversible compression. The decompression procedure delivers an estimated reproduction, which might be alluring where information of a few ranges that couldn't be perceived by the human mind can be dismissed. Such strategies could be utilized for media pictures, video and sound to accomplish more smaller information compression.

Lossless information compression is utilized to reduced records or information into a littler shape. It is frequently used to bundle up programming before it is sent over the Internet or downloaded from a site to decrease the measure of time and transmission capacity required to transmit the information. Lossless information compression has the imperative that when information is uncompressed, it must be indistinguishable to the first information that was compacted. Illustrations, sound, and video compression, for example, JPG, MP3, and MPEG then again utilize lossy compression plans which discard a portion of the first information to pack the documents considerably further. We will concentrate on the lossless kind. There are for the most part two classes of lossless compressors: word reference compressors and factual compressors. Word reference compressors, (for example, Lempel-Ziv based calculations) fabricate lexicons of strings and supplant whole gatherings of images. The factual compressors create models of the insights of the info information and utilize those models to control the last yield. Number juggling coding is a measurable compression system that utilizes evaluations of the probabilities of occasions to allocate code words. In a perfect world, short code words are relegated to more plausible occasions and longer code words are allotted to less likely occasions. Hypothetically, number-crunching codes allocate one "code word" to every

conceivable informational index. The number-crunching coder must cooperate with a modeler that gauges the probabilities of the occasions in the coding. To acquire great compression, a great likelihood show and productive method for speaking to the likelihood model are required. The models can be versatile, semi-versatile, or non-versatile. Versatile models progressively assess the presumably of every occasion in view of going before occasions. Semi-versatile models utilize a preparatory go of the information to accumulate a few measurements, and non-versatile models utilize settled probabilities for all information. Favorable position of number juggling coding is the partition of coding and demonstrating since it permits the multifaceted nature of the modeler to change without modifying the coder. The disservice is that it runs all the more gradually and is more intricate to execute than LZ based calculations.

There are three principle classes of lossless information compression methods: those utilizing measurable models, those that require the utilization of a word reference, and those that utilize both factual and lexicon based strategies. Word reference based compression plans have a tendency to be utilized more to archive applications (now and again in conjunction with different strategies), while continuous circumstances ordinarily require measurable compression plans. This is on the grounds that lexicon based calculations have a tendency to have moderate compression speeds and quick decompression speeds while measurable calculations have a tendency to be similarly quick amid compression and decompression. Factual compression plans decide the yield in light of the likelihood of event of the info images and are commonly utilized as a part of constant applications. The calculations have a tendency to be symmetric (the decoder reflects the means of the encoder); in this manner, compression and decompression for the most part require a similar measure of time to finish. Lexicon compression plans don't utilize a prescient measurable model to decide the likelihood of event of a specific image, yet they store strings of beforehand info images in a word reference. Word reference based compression methods are ordinarily utilized as a part of documenting applications, for example, pack and gzip on the grounds that the deciphering procedure has a tendency to be speedier than encoding. Cross breed compression strategies impart attributes to both factual and word reference based compression procedures. These calculations for the most part include a lexicon plot in a circumstance where rearranging suspicions can be made about the info information. Lossless information compression calculation incorporates: Lempel Ziv family (Dictionary-based encoding), Run-length Encoding compression (Statistical coding), Huffman Encoding (Statistical coding), Arithmetic Encoding (Statistical coding), Bitmask coding (Dictionary-based).

Number-crunching encoding is the most intense compression methods. The strategy for number juggling coding was recommended by Elms and exhibited by Abramson in his test on data hypothesis. In number-crunching coding a source group is introduced by an interim in the vicinity of 0 and 1 on the genuine number line. Useful executions of Arithmetic Coding are fundamentally the same as Huffman coding, despite the fact that it outperforms the Huffman system in its compression capacity. The Huffman strategy doles out an indispensable number of bits to every image, while number-crunching coding doles out one long code to the whole

info string. Number-crunching coding can possibly pack information to its hypothetical breaking point. Number juggling coding joins a measurable model with an encoding step, which comprises of a couple of math operations. The most fundamental factual model would have a straight time multifaceted nature of  $N[\log(n)+a] + Sn$  where N is the aggregate number of information images, n is the present number of special images, an is the math to be performed, and S is the time required, if vital, to keep up inward information structure. This changes over the whole information into a solitary skimming point number. A skimming guide number is comparative toward a number with a decimal point, as 4.5 rather than 41/2. In any case, in number-crunching coding we are not managing decimal number so we call it a skimming point rather than decimal point. The reason for information compression is the numerical estimation of data. Data contained in an image x is given by

$$L(x) = \text{Log}_2 \frac{1}{p(x)}$$

This value also describes the number of bits necessary to encode the symbol. This definition reinforces the notion of information. First, the more probable the occurrences of a symbol, the fewer bits are used to represent it. Conversely, the least frequent symbols provide more information by their occurrence. Secondly, if there are n equally probable messages,  $\log_2 n$  bits will be required to encode each message. This is the information value of each message

$$L = \text{Log}_2 \frac{1}{p(x)} = \log_2 n$$

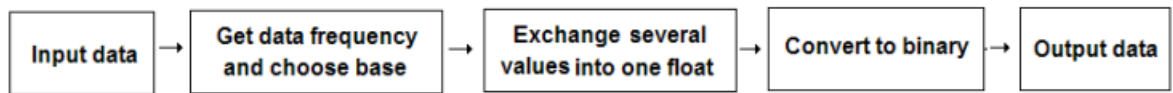


Fig 2-1: Arithmetic Data compression Flow Graph

New Character	Low value	High Value	Encoded Number	Output Symbol	Low	High	Range
	0.0	1.0	0.2572167752	B	0.2	0.3	0.1
B	0.2	0.3	0.572167752	I	0.5	0.6	0.1
I	0.25	0.26	0.72167752	L	0.6	0.8	0.2
L	0.256	0.258	0.6083876	L	0.6	0.8	0.2
L	0.2572	0.2576	0.041938	SPACE	0.0	0.1	0.1
SPACE	0.25720	0.25724	0.41938	G	0.4	0.5	0.1
G	0.257216	0.257220	0.1938	A	0.2	0.3	0.1
A	0.2572164	0.2572168	0.938	T	0.9	1.0	0.1
T	0.25721676	0.2572168	0.38	E	0.3	0.4	0.1
E	0.257216772	0.257216776	0.8	S	0.8	0.9	0.1
S	0.2572167752	0.2572167756	0.0				

(a)

(b)

Fig 2-2 (a) Encoding of “BILL GATES” (b) Decoding of “BILL GATES”

### 2.1.1 Practical Implementation

The way toward encoding and unraveling a surge of images utilizing number juggling coding is not very entangled. In any case, at first look, it appears to be totally unfeasible. Most PCs bolster skimming direct quantities of up toward 80 bits or something like that. Does this mean you need to begin once again every time you get done with encoding 10 or 15 images? Do you require a gliding point processor? Could machines with various gliding point designs impart utilizing number juggling coding?

For reasons unknown, number-crunching coding is best expert utilizing standard 16-bit and 32-bit whole number math. No gliding point math is required, nor would it utilize it. Rather, we utilize an incremental transmission conspire in which settled size whole number state factors get new bits at the low end and move them out the top of the line, framing a solitary number that can be the length of the quantity of bits accessible on the PC's stockpiling medium.

In the past area, I indicated how the calculation works by monitoring a high and low number that section the scope of the conceivable yield number. At the point when the calculation first begins up, the low number is set to 0.0, and the high number to 1.0. To work with number math, first change the 1.0 to 0.999..., or .111 ... in double.

To store these numbers in whole number registers, we first legitimize them so the suggested decimal point is on the left-hand side of the word. At that point we stack the same number of the underlying high and low values as will fit into the word estimate we are working with. My execution utilizes 16-bit unsigned math, so the underlying estimation of high is 0xFFFF, and low is 0. We realize that the high esteem proceeds with FFs perpetually, and low proceeds with 0s everlastingly, so we can move those additional bits in with exemption when they are required.

On the off chance that you envision our BILL GATES case in a 5-digit enlist, what might as well be called our setup would look like Figure 7(a). To locate our new range numbers, we have to apply the encoding calculation from the past area. We first figure the range between the low and high values. The contrast between the two registers will be 100000, not 99999, in light of the fact that expecting the high enroll has an unending number of 9s included to it, we have to increase the figured distinction. We then figure the new high esteem utilizing the recipe from the past segment:  $\text{high} = \text{low} + \text{high range (image)}$ .

For this situation the high range was .30, which gives another incentive for high of 30000. Before putting away the new estimation of high, we have to decrement it, by and by in light of the inferred digits added to the whole number esteem. So the new estimation of high is 29999. The count of low takes after a similar way, with a subsequent new estimation of 20000. So now high and low resemble this:

HIGH: 29999 (999...)

LOW: 20000 (000...)

Now, the most noteworthy digits of high and low match. Because of the way of our calculation, high and low can keep on growing more like each other without perpetually coordinating. This implies once they coordinate in the most critical digit, that digit will never show signs of change. So we can now yield that digit as the main digit of our encoded number. This is finished by moving both high and low left by one digit, and moving in a 9 at all huge digit of high. The comparable operations are performed in twofold in the C execution of this calculation. As this procedure proceeds with, high and low are consistently developing nearer together, then moving digits out into the coded word. The procedure for our "BILL GATES" message

	HIGH	LOW	RANGE	CUMULATIVE OUTPUT
Initial state	99999	00000	100000	
Encode B (0.2-0.3)	29999	20000		
Shift out 2	99999	00000	100000	.2
Encode I (0.5-0.6)	59999	50000		.2
Shift out 5	99999	00000	100000	.25
Encode L (0.6-0.8)	79999	60000	20000	.25
Encode L (0.6-0.8)	75999	72000		.25
Shift out 7	59999	20000	40000	.257
Encode SPACE (0.0-0.1)	23999	20000		.257
Shift out 2	39999	00000	40000	.2572
Encode G (0.4-0.5)	19999	16000		.2572
Shift out 1	99999	60000	40000	.25721
Encode A (0.1-0.2)	67999	64000		.25721
Shift out 6	79999	40000	40000	.257216
Encode T (0.9-1.0)	79999	76000		.257216
Shift out 7	99999	60000	40000	.2572167
Encode E (0.3-0.4)	75999	72000		.2572167
Shift out 7	59999	20000	40000	.25721677
Encode S (0.8-0.9)	55999	52000		.25721677
Shift out 5	59999	20000		.257216775
Shift out 2				.2572167752
Shift out 0				.25721677520

resembles this:

Note that after all the letters have been accounted for, two extra digits need to be shifted out of either the high or low value to finish up the output word.

### 2.1.2 Underflow

This plan functions admirably for incrementally encoding a message. There is sufficient exactness held amid the twofold accuracy whole number figurings to guarantee that the message is precisely encoded. Nonetheless, there is potential for lost exactness in specific situations.

If the encoded word has a series of 0s or 9s in it, the high and low values will gradually merge on an esteem, however may not see their most huge digits coordinate instantly. For instance, high and low may resemble this:

High: 700004  
Low: 699995

At this point, the calculated range is going to be only a single digit long, which means the output word will not have enough precision to be accurately encoded. Even worse, after a few more iterations, high and low could look like this :

High: 70000  
Low: 69999

Now, the qualities are forever stuck. The range amongst high and low has turned out to be small to the point that any count will dependably give back similar qualities. In any case, since the most critical digits of both words are not equivalent, the calculation can't yield the digit and move. It appears like an impasse.

The best approach to thrashing this undercurrent issue is to keep things from always getting this terrible. The first calculation said something like "If the most huge digit of high and low match, move it out". On the off chance that the two digits don't coordinate, yet are currently on adjoining numbers, a moment test should be connected. Assuming High and low are one separated, we then test to check whether the second most critical digit in high is a 0, and the second digit in low is a 9. Assuming this is the case, it implies we are headed for sub-current, and need to make a move.

At the point when undercurrent raises its appalling head, we take it off with a somewhat extraordinary move operation. Rather than moving the most critical digit out of the word, we simply erase the second digits from high and low, and move whatever is left of the digits left to top off the space. The most noteworthy digit remains set up. We then need to set an undercurrent counter to recall that we discarded a digit, and we aren't exactly certain whether it would wind up as a 0 or a 9. The operation resembles this:

	Before	After
High:	40344	43449
Low:	39810	38100
Underflow:	0	1

After each recalculation operation, if the most huge digits don't coordinate, we can check for sub-current digits once more. In the event that they are available, we move them out and increase the counter.



At the point when the most critical digits do at long last meet to a solitary esteem, we first yield that esteem. At that point, we yield the majority of the "sub-current" digits that were beforehand disposed of. The sub-current digits will be every one of the 9s or 0s, contingent upon whether High and Low focalized to the higher or lower esteem. In the C execution of this calculation, the undercurrent counter monitors what number of ones or zeros to put out.

COMPRESSION METHOD	ARITHMETIC	HUFFMAN
Compression ratio	Very good	Poor
Compression speed	Slow	Fast
Decompression speed	Slow	Fast
Memory space	Very low	Low
Compressed pattern matching	No	Yes
Permits Random access	No	Yes

Table 1: Comparison between arithmetic and Huffman coding methodologies

## 2.2 LZW

The accompanying case delineates the LZW calculation in real life, demonstrating the status of the yield and the word reference at each stage, both in encoding and unraveling the information. This case has been built to give sensible compression on a short message. In genuine content information, redundancy is for the most part less articulated, so longer information streams are regularly vital before the compression develops productivity.

The plaintext to be encoded (from a letters in order utilizing just the capital letters) is:

TOBEORNOTTOBEORTOBEORNOT#

The # is a marker used to demonstrate that the finish of the message has been come to. There are along these lines 26 images in the plaintext letter set (the 26 capital letters A through Z), in addition to the stop code #. We subjectively relegate these the qualities 1 through 26 for the letters, and 0 for '#'.

A PC will render these as series of bits. Five-piece codes are expected to give adequate blends to envelop this arrangement of 27 qualities. The word reference is introduced with these 27 values. As the word reference develops, the codes should develop in width to suit the extra passages. A 5-bit code gives  $2^5 = 32$  conceivable blends of bits, so when the 33rd lexicon word is made, the calculation should switch by then from 5-bit strings to 6-bit strings (for all code values, including those which were already yield with just five bits). Take note of that since the each of the zero code 00000 is utilized, and is named "0", the 33rd word reference section will be named 32. (Beforehand created yield is not influenced by the code-width change, but rather once a 6-bit esteem is produced in the word reference, it could possibly be the following code transmitted, so the width for resulting yield movements to 6 bits to suit that.)

The initial dictionary, then, will consist of the following entries:

Symbol	Binary	Decimal
#	00000	0
A	00001	1
B	00010	2
C	00011	3
D	00100	4
E	00101	5
F	00110	6
G	00111	7
H	01000	8
I	01001	9
J	01010	10
K	01011	11
L	01100	12
M	01101	13
N	01110	14
O	01111	15
P	10000	16

Q	10001	17
R	10010	18
S	10011	19
T	10100	20
U	10101	21
V	10110	22
W	10111	23
X	11000	24
Y	11001	25
Z	11010	26

### Encoding

Unencoded length = 25 symbols  $\times$  5 bits/symbol = 125 bits

Encoded length = (6 codes  $\times$  5 bits/code) + (11 codes  $\times$  6 bits/code) = 96 bits.

Utilizing LZW has spared 29 bits out of 125, decreasing the message by very nearly 22%. In the event that the message were longer, then the lexicon words would start to speak to longer and longer areas of content, permitting rehashed words to be sent minimalistically.

Current Sequence	Next Char	Output		Extended Dictionary		Comments
		Code	Bits			
NULL	T					
T	O	20	10100	27:	TO	27 = first available code after 0 through 26
O	B	15	01111	28:	OB	
B	E	2	00010	29:	BE	
E	O	5	00101	30:	EO	
O	R	15	01111	31:	OR	
R	N	18	10010	32:	RN	32 requires 6 bits, so for next output use 6 bits

N	O	14	001110	33:	NO	
O	T	15	001111	34:	OT	
T	T	20	010100	35:	TT	
TO	B	27	011011	36:	TOB	
BE	O	29	011101	37:	BEO	
OR	T	31	011111	38:	ORT	
TOB	E	36	100100	39:	TOBE	
EO	R	30	011110	40:	EOR	
RN	O	32	100000	41:	RNO	
OT	#	34	100010			# stops the algorithm; send the cur seq
		0	000000			and the stop code

### Decoding

To unravel a LZW-packed document, one has to know ahead of time the underlying word reference utilized, however extra sections can be remade as they are dependably just connections of past passages.

Input		Output Sequence	New Dictionary Entry				Comments
Bits	Code		Full		Conjecture		
10100	20	T			27:	T?	
01111	15	O	27:	TO	28:	O?	
00010	2	B	28:	OB	29:	B?	
00101	5	E	29:	BE	30:	E?	

01111	15	O	30:	EO	31:	O?	
10010	18	R	31:	OR	32:	R?	created code 31 (last to fit in 5 bits)
001110	14	N	32:	RN	33:	N?	so start reading input at 6 bits

## 2.3

## Run Length Encoding

Run-length encoding is an information compression calculation that is upheld by most bitmap document arrangements, for example, TIFF, BMP, and PCX. RLE is suited for compacting any sort of information paying little respect to its data content, yet the substance of the information will influence the compression proportion accomplished by RLE. Albeit most RLE calculations can't accomplish the high compression proportions of the more propelled compression techniques, RLE is both simple to actualize and brisk to execute, making it a decent other option to either utilizing a perplexing compression calculation or leaving your picture information uncompressed.

RLE works by decreasing the physical size of a rehashing series of characters. This rehashing string, called a run, is ordinarily encoded into two bytes. The primary byte speaks to the quantity of characters in the run and is known as the run number. By and by, an encoded run may contain 1 to 128 or 256 characters; the run consider more often than not contains the quantity of characters short one (an incentive in the scope of 0 to 127 or 255). The second byte is the estimation of the character in the run, which is in the scope of 0 to 255, and is known as the run esteem.

Uncompressed, a character keep running of 15 A characters would regularly require 15 bytes to store:

AAAAAAAAAAAAAAAA

A similar string after RLE encoding would require just two bytes:

15A

The 15A code created to speak to the character string is called a RLE bundle. Here, the main byte, 15, is the run check and contains the quantity of reiterations. The second byte, An, is the run esteem and contains the real rehashed an incentive in the run.

Another bundle is created each time the run character changes, or each time the quantity of characters in the run surpasses the most extreme number. Expect that our 15-character string now contains four distinctive character runs:

```
AAAAAAbbbXXXXXt
```

Utilizing run-length encoding this could be compacted into four 2-byte bundles:

```
6A3b5X1t
```

In this way, after run-length encoding, the 15-byte string would require just eight bytes of information to speak to the string, instead of the first 15 bytes. For this situation, run-length encoding yielded a compression proportion of very nearly 2 to 1.

Long runs are uncommon in specific sorts of information. For instance, ASCII plaintext from time to time contains long runs. In the past illustration, the last run (containing the character t) was just a solitary character long; a 1-character run is as yet a run. Both a run check and a run esteem must be composed for each 2-character run. To encode a keep running in RLE requires at least two characters worth of data; in this way, a keep running of single characters really consumes more room. For similar reasons, information comprising completely of 2-character runs continues as before size after RLE encoding.

In our illustration, encoding the single character toward the end as two bytes did not observably hurt our compression proportion in light of the fact that there were such a variety of long character keeps running in whatever is left of the information. Be that as it may, watch how RLE encoding duplicates the span of the accompanying 14-character string:

```
Xtmprdqzntwlfb
```

After RLE encoding, this string progresses toward becoming:

```
1X1t1m1p1r1d1q1z1n1t1w1l1f1b
```

RLE plans are basic and quick, however their compression effectiveness relies on upon the sort of picture information being encoded. A highly contrasting picture that is generally white, for example, the page of a book, will encode extremely well, because of the huge measure of coterminous information that is all a similar shading. A picture with many hues that is exceptionally occupied in appearance, in any case, for example, a photo, won't encode extremely well. This is on account of the unpredictability of the picture is communicated as countless hues. Also, in view of this unpredictability there will be moderately few keeps running of a similar shading.

## 2.4 Deflate

The information is compacted as a blend of encoded bytes and coordinating strings, where the strings are to be found in the first uncompressed information. Each match is a length and a separation over from the present position. The literals and lengths are consolidated into a solitary Huffman code, and the separations in another Huffman code. Longer lengths and separations fall into containers, trailed by additional bits to figure out which section in the canister to utilize. The stream comprises of a progression of exacting/length codes, where a length code is trailed by a separation code. A separation might be not as much as the length, in which case the past accessible information is replicated, and after that what was duplicated is replicated again until the length is come to. The lengths can be in 3.258, and the separations can be in 1.32768, where 32768 bytes is the measure of past information held. This general way to deal with code a succession of literals and matches is called "LZ77".

The flatten stream is broken into pieces, where each square begins with the meaning of the Huffman codes for that square, trailed by the strict/length and separation codes, lastly an end-of-piece code (a unique exacting/length code). The depiction of the Huffman codes comprises of the code lengths of every image, where that portrayal is itself Huffman and run-length encoded. The last piece is set apart in that capacity, so the empty configuration is self-ending.

There are additionally settled Huffman code squares, which utilize a solitary pre-characterized strict/length and separation code, and put away pieces which basically duplicate the information uncompressed.

The collapse calculation scans for coordinating strings in the previous information. It is in this pursuit that the greater part of the viability of the compressor is resolved, and also its speed. In zlib a hash table is built on all former three-byte sequences (clearing old ones as the window slides out). Hash hits result in straightforwardly searching for a coordinating string at the related separation back. The culmination of the hash table and when to quit searching for better matches, and thusly the speed, is controlled by the client chose compression level. Other flatten compressors utilize more entire, yet slower methodologies, for example, postfix trees to discover past matches.

The collapse viability can likewise be upgraded by intelligently picking when to begin another piece, which permits the calculation to adjust to changing measurements in the information.

The Huffman codes are produced utilizing, actually, the Huffman calculation, however changed to be compelled to a length farthest point of 15 bits.

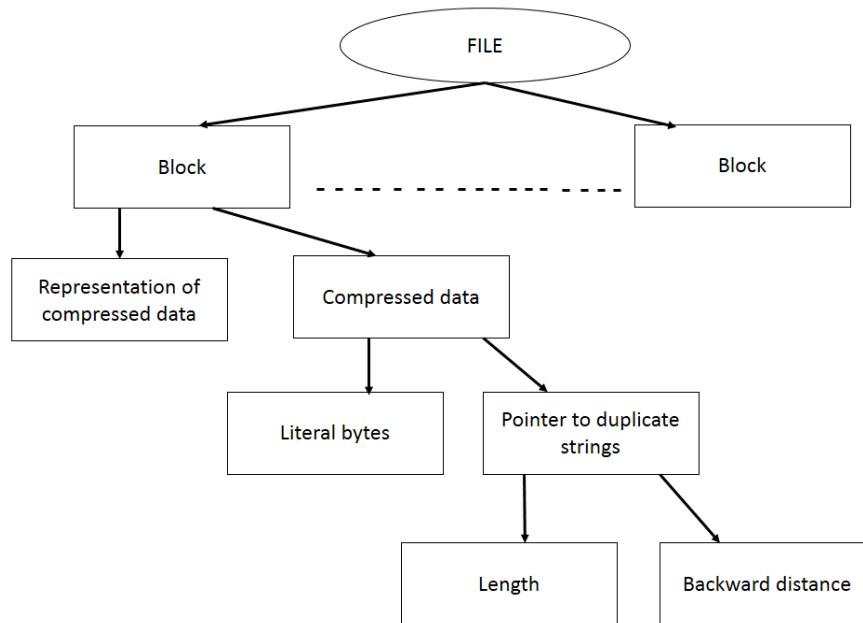


Fig how the file will be processed into streams.

## 2.5 Prediction by partial matching

PPM was created in 1984 by Cleary and Witten not long after math encoding was proficient in 1976. This was trailed by a progression of changes (Moffat, 1990) bringing about the variant PPMC. The calculation utilizes a limited setting Markov Model to anticipate the event of the following image in view of the events of the past images in the content used to prepare the model. As a general rule, the technique really mixes together numerous lower arrange Markov models for the situation when a given request is inadmissible for the forecast. In this way it recursively falls onto bring down request models until the slightest request where each image happens with a likelihood of  $1/|\Sigma|$  where  $\Sigma$  is the aggregate number of images in the letters in order. This records for the escape likelihood on the off chance that an image was at no other time 20 found in the information content used to create the model insights. At its presentation, PPM was thought to have the capacity to join any request demonstrate and the forecasts would enhance with the expanding request. In any case it was later recognized that a request of "5" works best and ideally and as the request increments facilitate, the exactness falls attributable to the escape component onto bring down request models. The PPM calculation that uses an unbound setting length shows it is ideal to utilize limited deterministic settings for application such grouping of content (Hiroyuki et al, 2005). In spite of the fact that the underlying utilizations of the PPM calculation were to empower better compression, it has been connected



to content forecast, dialect distinguishing proof, tongue recognizable proof, dialect division, word division, content classification and so on. Give us a chance to consider an information content "abracadabra". We might now watch how PPM will utilize the image frequencies to produce models of request  $k=2, 1, 0,$  and  $-1$  and utilize these disseminations to foresee the likelihood of event of the following image in the unique situation.

Order $k = 2$			Order $k = 1$			Order $k = 0$			Order $k = -1$		
Prediction	c	p	Prediction	c	p	Prediction	c	p	Prediction	c	p
ab →r →esc	2 1	2/3 1/3	a →b →c →d →esc	2 1 1 3	2/7 1/7 1/7 3/7	→a	5	5/16	→A	1	1/ A
ac →a →esc	1 1	1/2 1/2	b →r →esc	2 1	2/3 1/3	→b	2	2/16			
ad →a →esc	1 1	1/2 1/2	c →a →esc	1 1	1/2 1/2	→c	1	1/16			
ba →a →esc	2 1	2/3 1/3	d →a →esc	1 1	1/2 1/2	→d	1	1/16			
ca →d →esc	1 1	1/2 1/2	r →a →esc	2 1	2/3 1/3	→r	2	2/16			
da →b →esc	1 1	1/2 1/2				→esc	5	5/16			
ra →c →esc	1 1	1/2 1/2									

## Chapter-3

### SYSTEM DEVELOPMENT

#### 3.1 Arithmetic Coding

##### 3.1.1 Compression algorithm

```

begin
  count frequency of input symbols
  output (symbol's frequency)
  interval I := new interval 0..9999
  split I according frequency of symbols
  readSymbol(X)
  while (X!=EOF) do
    begin
      if (MSB(Low) == MSB(High) then
        begin
          output(MSB)
          in case of need output discarded digits
          shift left High and Low
        end
      else
        begin
          if (underflow danger) then
            begin
              shift left High and Low from second position
            end
          end
          new I := interval I accordant with X
          split I according frequency of symbols
          readSymbol(X)
        end
      end
      output(remainder)
    end
  end
end

```

### 3.1.2 Decompression algorithm

- $\text{Range} = (\text{high} - \text{low}) + 1$  See where the number lands
- $\text{Temp} = ((\text{code} - \text{low}) + 1) * \text{scale} - 1) / \text{range}$
- See what symbols corresponds to temp.
- $\text{Range} = (\text{high} - \text{low}) + 1$  Extract the symbol code
- $\text{High} = \text{low} + ((\text{range} * \text{high\_values}[\text{symbol}]) / \text{scale}) - 1$
- $\text{Low} = \text{low} + (\text{range} * \text{high\_values}[\text{symbol} - 1]) / \text{scale}$  Note that those formulae are the same that the encoder uses
- Loop.
- Msb of high = msb of low?
- Yes
  - Go to shift
- No
  - Second msb of low = 1 and Second msb of high = 0 ?
  - Yes
    - $\text{Code} = \text{code} \wedge 4000\text{h}$
    - $\text{Low} = \text{low} \& 3\text{FFFh}$
    - $\text{High} = \text{high} | 4000\text{h}$
    - go to shift
  - No
    - The routine for decoding a symbol ends here.

#### Shift:

- Shift low to the left one time. Now we have to put in low, high and code new bits
- Shift high to the left one time, and or the lsb with the value 1
- Shift code to the left one time, and or it the next bit in the input
- Repeat to the first loop.

### 3.1.3 Model development

The need to precisely foresee the likelihood of images in the information is characteristic to the way of number juggling coding. The standard of this sort of coding is to lessen the quantity of bits expected to encode a character as its likelihood of appearance increments. So if the letter "e" speaks to 25 percent of the information, it would just take 2 bits to code. On the off chance that the letter "z" speaks to just 0.1 percent of the info information, it may take 10 bits to code. In the event that the model is not producing probabilities precisely, it may take 10 bits to speak to "e" and 2 bits to speak to "z," bringing on information extension rather than compression.

The second condition is that the model needs to make forecasts that go amiss from a uniform appropriation. The better the model is at making these expectations, the better the compression proportions will be. For instance, a model could be made that allocated every one of the 256 conceivable images a uniform likelihood of  $1/256$ . This model would make a yield document that was the very same size as the information record, in light of the fact that each image would take precisely 8 bits to encode. Just by accurately discovering probabilities that stray from a

uniform appropriation can the quantity of bits be decreased, prompting compression. Obviously, the expanded probabilities need to precisely reflect reality, as recommended by the principal condition.

The most proficient system for processing dispersions relies on upon the information sort. When we are managing totally obscure information we may need adjustment to work in a totally programmed way. In different cases, we can utilize some learning of the information properties to diminish or dispose of the adjustment exertion. Beneath we clarify the components of the absolute most basic methodologies for evaluating dispersions.

- Use a steady dissemination that is accessible before encoding and interpreting, ordinarily assessed by social occasion insights in countless examples. This approach can be utilized for sources, for example, English content, or climate information, yet it once in a while yields the best outcomes since few data sources are so basic as to be demonstrated by a solitary circulation. Moreover, there is next to no adaptability (e.g., insights for English content don't fit well Spanish content). Then again, it might function admirably if the source model is extremely point by point, and in actuality it is the main option in some exceptionally complex models in which important measurements must be assembled from a lot of information.
- Use pre-characterized circulations with versatile parameter estimation. For example, we can expect that the information has Gaussian circulation, and gauge just the mean and difference of every image. In the event that we permit just a couple values for the circulation parameters, then the encoder and decoder can make a few vectors with all the dispersion values, and utilize them as indicated by their basic parameter estimation.
- Use two-pass encoding. A first pass assembles the insights of the source, and the second pass codes the information with the gathered measurements. For interpreting, a scaled rendition of vectors  $p$  or  $c$  must be incorporated toward the start of the packed information. For instance, a book can be filed (compacted) together with its specific image insights. It is conceivable to lessen the computational overhead by sharing procedures between passes. For instance, the primary pass can at the same time assemble insights and change over the information to run-lengths.
- Use a dissemination in view of the event of images beforehand coded, refreshing  $c$  with every image encoded. We can begin with an extremely inexact circulation (e.g., uniform), and if the probabilities change as often as possible, we can reset the appraisals occasionally. This system, clarified in the following segment, is very successful and the most helpful and adaptable. Be that as it may, the steady refresh of the aggregate dissemination can expand the computational multifaceted nature significantly. An option is to refresh just the likelihood vector  $p$  after each encoded image, and refresh the combined appropriation  $c$  less as often as possible.

#### 3.1.4 Analysis of arithmetic coding

Number-crunching coding is notable for its optimality, and the way that it can be an exceptionally flexible and intense apparatus for coding complex information sources [1, 2, 4, 6, 10]. In the meantime, specialists likewise realize that it had not been all the more regularly utilized due to its high computational multifaceted nature. On the off chance that we consider the numerous times of research on procedures for diminishing its multifaceted nature, it might appear that there is little seek after new achievements and for its boundless appropriation. Notwithstanding, new outcomes demonstrate that, truth be told, the development of number-crunching coding is taking after a surprising way, and the most encouraging option is to move back to the easiest executions. This happens on the grounds that the vast majority of the cost-diminishment systems for number juggling coding were produced for the equipment that was accessible at least 10 years back, when augmentations and divisions were too moderate for coding purposes. At present even reasonable processors can perform exact number juggling quick. Since number juggling is such a key errand of any processor, we can expect considerably more prominent points of interest later on. In this way, there is a need to distinguish what are the number juggling coding errands that will remain really critical in deciding the computational multifaceted nature. Utilizing this data we ought to have the capacity to discover how to better adventure the processor's math abilities for quicker coding. The wellsprings of number juggling coding computational unpredictability incorporate [1, 2]:

- Interval refresh and math operations
- Symbol deciphering (interim hunt)
- Interval renormalization
- Carry spread and bit moves
- Probability estimation (source demonstrating)

Since every one of these activities can be firmly incorporated in a solitary usage, it has been hard to unmistakably recognize the execution bottlenecks. In this work we handle this issue by doing a broad relative assessment. Our methodology is to gauge the execution of a few usage, transforming one parameter at time, or if nothing else as few as could reasonably be expected. Along these lines we can assess the significance of each of the errands specified above, and furthermore select the best procedures.

While the principle goal of this paper is to assess the distinction in intricacy of an assortment of errands and strategies, in a few diagrams we include comes about because of some notable executions, since they can give a flat out reference (we clarify when the correlations are not reasonable). To keep away from repetition, in this report we don't present every one of the points of interest of our execution and examination, since a large portion of it can be found in references [1, 2] (which likewise gives an a great deal more entire introduction, and set of

references on number-crunching coding hypothesis and practice). Notwithstanding, the peruser ought to know that there is a lot of programming for each investigation: we needed to compose particular projects to test all the essential assignments, in more than 10 diverse full executions of number juggling coding. Despite the fact that a portion of the procedures for many-sided quality lessening we display here are not new, we trust this is the first occasion when that their utilization for number-crunching coding is accounted for in this type of unpredictability examinations. Furthermore, we trust that the consecutive detachment of the examination of the diverse undertakings, with the distinguishing proof of the most suited usage to be utilized as a part of different tests empowered a greatly improved comprehension of the many-sided quality issues.

Math coders deliver close ideal yield for a given arrangement of images and probabilities. Compression calculations that utilization number juggling coding begin by deciding a model of the information - essentially an expectation of what examples will be found in the images of the message. The more exact this expectation is, the nearer to optimality the yield will be.

Case: a straightforward, static model of information:

- 60% possibility of image "a" - > the interim would be [0, 0.6)
- 20% possibility of image "b" - > the interim would be [0.6, 0.8)
- 10% possibility of image "c" - > the interim would be [0.8, 0.9)
- 10% possibility of image END-OF-DATA. - > the interim would be [0.9, 1)

The nearness of END-OF-DATA image implies that the stream will be 'inside ended', as is genuinely regular in information compression; the first and final time this image shows up in the information stream, the decoder will realize that the whole stream has been decoded.

The encoder has fundamentally only three bits of information to consider: the following image that should be encoded, the present interim, the probabilities of images. Because of this, is very simple to alter the calculation to versatile model.

The encoder isolates the present interim into sub-interims, each speaking to a small amount of the present interim corresponding to the likelihood of that image in the present setting. Whichever interim relating to the genuine image that is by be encoded turns into the interim utilized as a part of the following step. When the sum total of what images have been encoded, the subsequent interim distinguishes, unambiguously, the arrangement of images that created it. Any individual who has the last interim and the model utilized can remake the image arrangement that more likely than not been entered the encoder to bring about that last interim. Memory complexity depends on a number of different input symbols, at maximum

$O(n)$ , where  $n$  is length of a message. Time complexity depends on a number of different input symbols and length of a message. So  $n + n*|\Sigma|$ , where  $|\Sigma|$  is a number of different input symbols, at maximum  $O(n*|\Sigma|)$ .

### 3.2 Dictionary-based Algorithm

#### 3.2.1 Compression Algorithm

1. STRING = get input character
2. WHILE not end of input stream DO
3.     CHARACTER = get input character
4.     IF STRING+CHARACTER is in the string table then
5.         STRING = STRING+CHARACTER
6.     ELSE
7.         Output the code for STRING
8.         add STRING+CHARACTER to the STRING table
9.         STRING = CHARACTER
10.    END of IF
11. END of WHILE
12. Output the code for STRING

#### 3.2.2 Decompression Algorithm

1. Read OLD\_CODE
2. output OLD\_CODE
3. CHARACTER = OLD\_CODE
4. WHILE there are still input characters DO
5.     Read NEW\_CODE
6.     IF NEW\_CODE is not in the translation table THEN
7.         STRING = get translation of OLD\_CODE
8.         STRING = STRING+CHARACTER
9.     ELSE
10.         STRING = get translation of NEW\_CODE
11.    END of IF
12. output STRING
13. CHARACTER = first character in STRING
14. add OLD\_CODE + CHARACTER to the translation table
15. OLD\_CODE = NEW\_CODE
16. END of WHILE

#### 3.2.3 Model development

The situation portrayed by Welch's 1984 paper[22] encodes groupings of 8-bit information as settled length 12-bit codes. The codes from 0 to 255 speak to 1-character successions comprising of the comparing 8-bit character, and the codes 256 through 4095 are made in a word reference for arrangements experienced in the information as it is encoded. At each phase in compression, input bytes are assembled into an arrangement until the following character would make a succession for which there is no code yet in the word reference. The code for the succession (without that character) is added to the yield, and another code (for the grouping with that character) is added to the word reference.

The thought was immediately adjusted to different circumstances. In a picture in light of a shading table, for instance, the characteristic character letter set is the arrangement of shading table records, and in the 1980s, many pictures had little shading tables (on the request of 16 hues). For such a diminished letters in order, the full 12-bit codes yielded poor compression unless the picture was vast, so the possibility of a variable-width code was presented: codes regularly begin one piece more extensive than the images being encoded, and as each code size is spent, the code width increments by 1 bit, up to some endorsed greatest (commonly 12 bits). At the point when the most extreme code esteem is achieved, encoding continues utilizing the current table, however new codes are not produced for expansion to the table.

Advance refinements incorporate saving a code to demonstrate that the code table ought to be cleared and reestablished to its underlying state (an "unmistakable code", normally the main esteem instantly after the qualities for the individual letter set characters), and a code to show the finish of information (a "stop code", regularly one more noteworthy than the reasonable code). The unmistakable code permits the table to be reinitialized after it tops off, which gives the encoding a chance to adjust to changing examples in the information. Brilliant encoders can screen the compression productivity and gather the dishes at whatever point the current table no longer matches the information well.

Since the codes are included a way controlled by the information, the decoder imitates building the table as it sees the subsequent codes. It is important that the encoder and decoder concede to which assortment of LZW is being utilized: the measure of the letters in order, the most extreme table size (and code width), regardless of whether variable-width encoding is being utilized, the underlying code estimate, whether to utilize the unmistakable and stop codes (and what values they have). Most configurations that utilize LZW incorporate this data with the organization determination or give express fields to them in a compression header for the information.

## Encoding

An abnormal state perspective of the encoding calculation is appeared here:

- Initialize the word reference to contain all strings of length one.



- Find the longest string  $W$  in the word reference that matches the present information.
- Emit the lexicon list for  $W$  to yield and expel  $W$  from the info.
- Add  $W$  took after by the following image in the contribution to the lexicon.
- Go to Step 2.

A lexicon is introduced to contain the single-character strings relating to all the conceivable info characters. The calculation works by looking over the information string for progressively longer substrings until it discovers one that is not in the word reference. At the point when such a string is found, the file for the string without the last character is recovered from the word reference and sent to yield, and the new string is added to the lexicon with the following accessible code. The last information character is then utilized as the following beginning stage to check for substrings.

Along these lines, progressively longer strings are enlisted in the lexicon and made accessible for ensuing encoding as single yield qualities. The calculation works best on information with rehashed designs, so the underlying parts of a message will see little compression. As the message develops, in any case, the compression proportion tends asymptotically to the maximum.[22].

## Decoding

The disentangling calculation works by perusing an incentive from the encoded input and yielding the relating string from the instated lexicon. Keeping in mind the end goal to remake the word reference in an indistinguishable route from it was worked amid encoding, it likewise acquires the following an incentive from the info and adds to the lexicon the connection of the present string and the main character of the string got by disentangling the following information esteem, or the principal character of the string simply yield if the following worth can't be decoded. The decoder then continues to the following information esteem and rehashes the procedure until there is no more contribution, and soon thereafter the last information esteem is decoded with no more increments to the word reference. Along these lines the decoder develops a word reference which is indistinguishable to that utilized by the encoder, and utilizations it to interpret consequent information values. In this way the full lexicon does not need be sent with the encoded information; simply the underlying word reference containing the single-character strings is adequate.

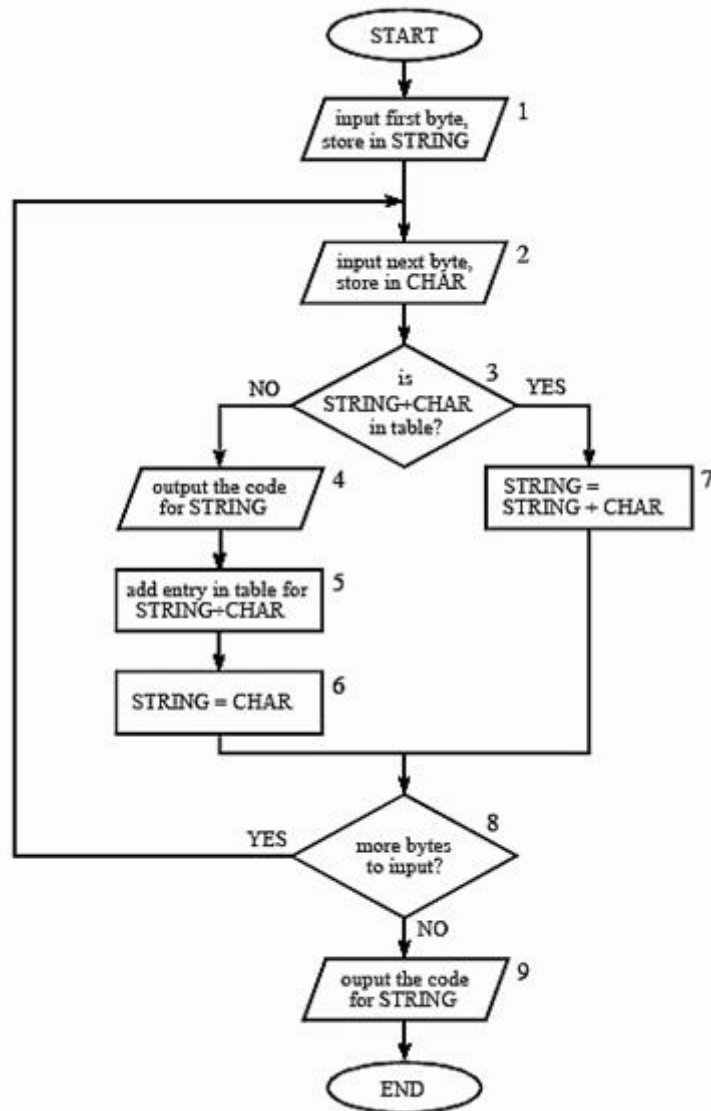


Fig 3-1: Flow chart of compression Algorithm

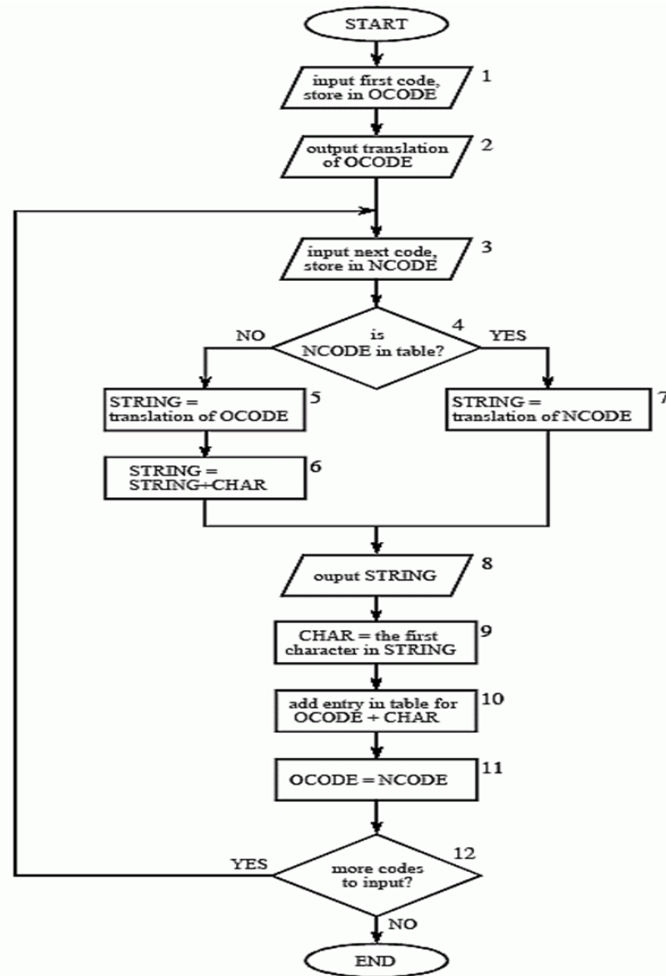


Fig 3-2: Flow chart of decompression Algorithm

### 3.2.4 Analysis of LZW

In 1984, Terry Welch was chipping away at a compression calculation for elite plate controllers. He built up a fairly basic calculation that depended on the LZ78 calculation and that is presently called LZW.

LZW compression replaces series of characters with single codes. It doesn't do any examination of the approaching content. Rather, it just includes each new series of characters it sees to a table of strings. Compression happens when a solitary code is yield rather than a series of characters. The

code that the LZW calculation yields can be of any subjective length, however it must have a larger number of bits in it than a solitary character. The initial 256 codes (when utilizing eight

piece characters) are naturally appointed to the standard character set. The rest of the codes are relegated to strings as the calculation continues. The example program keeps running as appeared with 12 bit codes. This implies codes 0-255 allude to individual bytes, while codes 256-4095 allude to substrings.

Favorable circumstances and disservices

- LZW compression works best for documents containing bunches of dreary information. This is regularly the case with content and monochrome pictures. Documents that are compacted however that don't contain any dull data whatsoever can even become greater!
- LZW compression is quick.
- LZW is a genuinely old compression system. All current PC frameworks have the pull to utilize more productive calculations.
- Royalties must be paid to utilize LZW compression calculations inside applications (see underneath).

LZW compression turned into the main broadly utilized all inclusive information compression technique on PCs. An expansive English content document can normally be compacted by means of LZW to about a large portion of its unique size.

LZW was utilized as a part of the general population area program pack, which turned into a pretty much standard utility in Unix frameworks around 1986. It has since vanished from numerous appropriations, both in light of the fact that it encroached the LZW patent and in light of the fact that gzip delivered better compression proportions utilizing the LZ77-based DEFLATE calculation, however starting at 2008 at any rate FreeBSD incorporates both pack and uncompress as a piece of the circulation. A few other well known compression utilities additionally utilized LZW, or firmly related strategies.

LZW turned out to be broadly utilized when it turned out to be a piece of the GIF picture design in 1987. It might likewise (alternatively) be utilized as a part of TIFF and PDF records. (In spite of the fact that LZW is accessible in Adobe Acrobat programming, Acrobat as a matter of course uses DEFLATE for most content and shading table-based picture information in PDF records.)

As the word reference size is settled and free of the info length, LZW is in  $O(n)$  as every byte is just perused once and the multifaceted nature of the operation for each character is steady. The LZW compression lexicon has tree character. On the off chance that put away in like manner, the word reference can be navigated one hub for every information byte, basically making the compression calculation  $O(n)$ - time in light of info length. Putting away the lexicon

that way most likely squanders heaps of memory, so it's the typical speed-space exchange off and a memory-productive execution likely is at any rate  $O(n \log n)$ .

### 3.1 Run Length Encoding

#### 3.1.1 Compression algorithm

```
// for input b with length l
// for input b with length l
index = 0
while( index < l )
{
    run = b[index]
    length = 0
    if run = b[++index]
        while run = b[index+length]
            index++, length++;
    output (run, length);
}
```

#### 3.1.2 Decompression algorithm

```
// for input b with length l
index = 0
while( index < l )
{
    run = b[index++]
    length = b[index++]
    output (run, length+1)
}
```

#### 3.1.3 Model Development

Model development of this algorithm is shown in the form of flow chart.

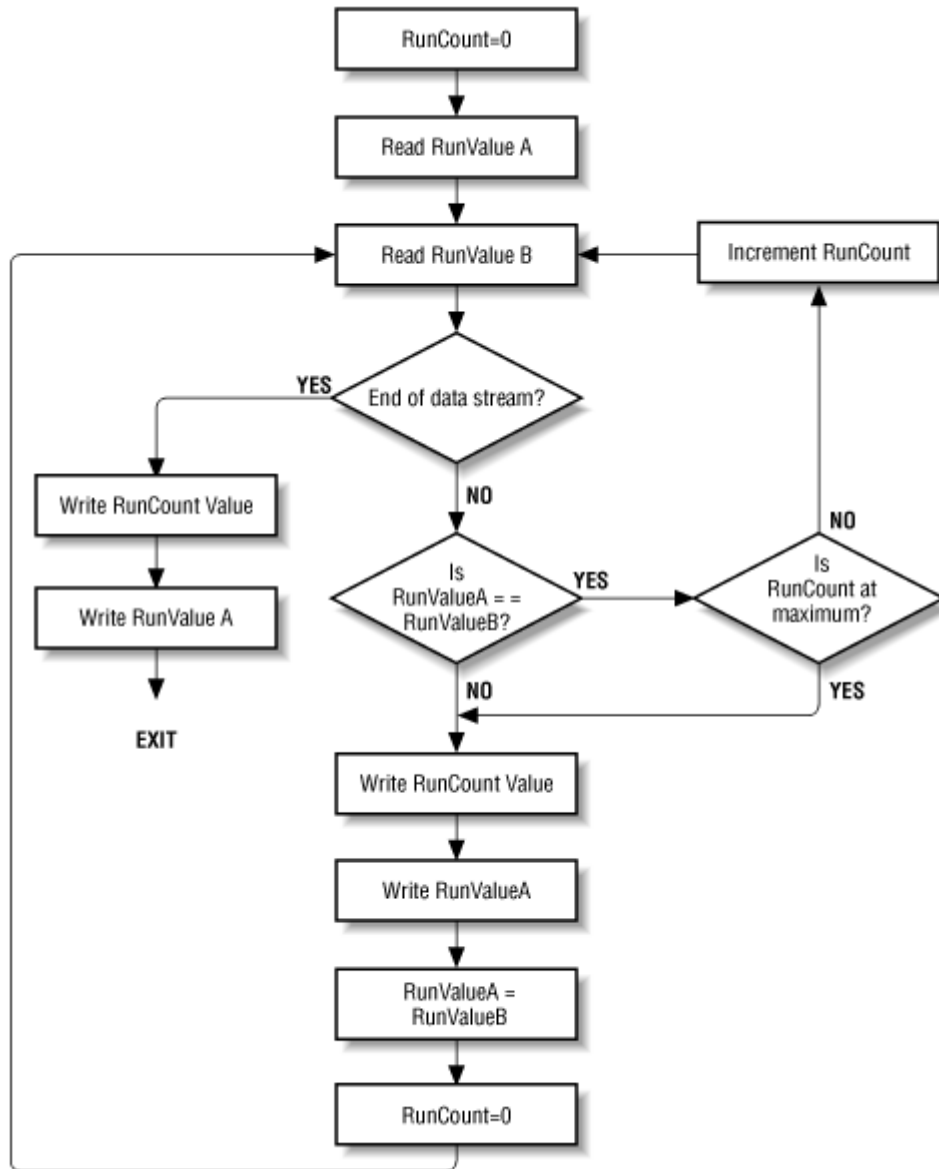


Fig 3-3: Basic flow chart of Run length Encoding Algorithm

### 3.1.4 Analysis of run length encoding

This calculation is anything but difficult to execute and does not require much CPU pull. RLE compression is just proficient with records that contain heaps of redundant information. These can be content documents on the off chance that they contain bunches of spaces for indenting yet line-workmanship pictures that contain vast white or dark zones are much more

appropriate. PC created shading pictures (e.g. building drawings) can likewise give reasonable compression proportions.

The parts of run-length encoding calculations that contrast are the choices that are made in light of the sort of information being decoded, (for example, the length of information runs). RLE plans used to encode bitmap design are normally separated into classes by the sort of nuclear (that is, most principal) components that they encode. The three classes utilized by most design document configurations are bit-, byte-, and pixel-level RLE.

Bit-level RLE plans encode keeps running of various bits in an output line and overlook byte and word limits. Just monochrome (high contrast), 1-bit pictures contain an adequate number of bit hurries to make this class of RLE encoding productive. An average piece level RLE plot encodes keeps running of one to 128 bits long in a solitary byte bundle. The seven slightest noteworthy bits contain the run tally less one, and the most huge piece contains the estimation of the bit run, either 0 or 1. A run longer than 128 pixels is part over a few RLE-encoded bundles.

Byte-level RLE plans encode keeps running of indistinguishable byte values, overlooking individual bits and word limits inside a sweep line. The most well-known byte-level RLE conspire encodes keeps running of bytes into 2-byte parcels. The principal byte contains the run check of 0 to 255, and the second byte contains the estimation of the byte run. It is additionally normal to supplement the 2-byte encoding plan with the capacity to store exacting, unencoded keeps running of bytes inside the encoded information stream too.

In such a plan, the seven minimum huge bits of the main byte hold the run check less one, and the most huge piece of the primary byte is the pointer of the sort of run that takes after the run number byte. On the off chance that the most noteworthy piece is set to 1, it means an encoded run. Encoded runs are decoded by perusing the run esteem and rehashing it the quantity of times showed by the run tally. On the off chance that the most noteworthy piece is set to 0, an exacting run is shown, implying that the following run check bytes are perused truly from the encoded picture information. The run number byte then holds an incentive in the scope of 0 to 127 (the run check short one). Byte-level RLE plans are useful for picture information that is put away as one byte for each pixel.

Pixel-level RLE plans are utilized when at least two back to back bytes of picture information are utilized to store single pixel qualities. At the pixel level, bits are disregarded, and bytes are numbered just to recognize every pixel esteem. Encoded parcel sizes change contingent on the span of the pixel qualities being encoded. The quantity of bits or bytes per pixel is put away in the picture document header. A keep running of picture information put away as 3-byte pixel values encodes to a 4-byte bundle, with one run-check byte took after by three run-esteem bytes. The encoding strategy continues as before as with the byte-situated RLE.

### 3.1 Deflate

#### 3.1.1 Compression algorithm

- Use of Huffman Coding
  - Shorter codes lexicographically precede longer codes.

Assuming that the order of the alphabet is ABCD:

Symbol	Code
A	10
B	0
C	110
D	111

i.e., 0 precedes 10 which precedes 11x, and 110 and 111 are lexicographically consecutive.

- 1) Count the number of codes for each code length. Let `bl_count[N]` be the number of codes of length `N`,  $N \geq 1$ .
- 2) Find the numerical value of the smallest code for each code length:

```

code = 0;
bl_count[0] = 0;
for (bits = 1; bits <= MAX_BITS; bits++) {
    code = (code + bl_count[bits-1]) << 1;
    next_code[bits] = code;
}

```

- 3) Assign numerical values to all codes, using consecutive values for all codes of the same length with the base values determined at step 2. Codes that are never used (which have a bit length of zero) must not be assigned a value.

```

for (n = 0; n <= max_code; n++) {
    len = tree[n].Len;
    if (len != 0) {
        tree[n].Code = next_code[len];
        next_code[len]++; } }

```



- Use of LZ77

```

begin
  fill view from input
  while (view is not empty) do
    begin
      find longest prefix p of view starting in coded part
      i := position of p in window
      j := length of p
      X := first char after p in view
      output(i,j,X)
      add j+1 chars
    end
  end
end

```

### 3.1.2 Decompression algorithm

```

do
  read block header from input stream.
  if stored with no compression
    skip any remaining bits in current partially processed byte
    copy bytes of data to output
  otherwise
    loop (until end of block code recognized)
      decode literal/length value from input stream
      if value < 256
        copy value (literal byte) to output stream
      otherwise
        if value = end of block (256)
          break from loop
        otherwise (value = 257..285)
          decode distance from input stream
          move backwards distance bytes in the output
          stream, and copy length bytes from this
          position to the output stream.
        end loop
    while not last block

```

### 3.1.3 Model Development

A Deflate stream comprises of a progression of pieces. Each piece is gone before by a 3-bit header:

- First bit: Last-obstruct in-stream marker:
- 1: this is the last piece in the stream.
- 0: there are more pieces to handle after this one.
- Second and third bits: Encoding strategy utilized for this piece sort:
- 00: a put away/crude/exacting segment, in the vicinity of 0 and 65,535 bytes long.
- 01: a static Huffman compacted piece, utilizing a pre-concurred Huffman tree.
- 10: a compacted piece finish with the Huffman table provided.
- 11: held, don't utilize.

The put away square choice includes insignificant overhead, and is utilized for information that is incompressible.

Most compressible information will wind up being encoded utilizing technique 10, the element Huffman encoding, which creates an upgraded Huffman tree altered for each square of information independently. Guidelines to produce the important Huffman tree instantly take after the square header. The static Huffman choice is utilized for short messages, where the settled sparing picked up by discarding the tree exceeds the rate compression misfortune because of utilizing a non-ideal (in this way, not in fact Huffman) code.

Compression is accomplished through two stages:

- The coordinating and supplanting of copy strings with pointers.
- Replacing images with new, weighted images in view of recurrence of utilization.

### 3.1.4 Analysis of deflate

While it is the plan of this archive to characterize the "flatten" packed information design without reference to a specific compression calculation, the arrangement is identified with the compacted groups delivered by LZ77.

The compressor ends a piece when it establishes that beginning a square with new trees would be helpful, or when the piece scrutinize fills the compressor's piece cushion.

The compressor utilizes an affixed hash table to discover copied strings, utilizing a hash capacity that works on 3-byte successions. At any given point amid compression, let XYZ be

the following 3 input bytes to be inspected (not really all extraordinary, obviously). To start with, the compressor inspects the hash chain for XYZ. On the off chance that the chain is unfilled, compressor just works out X as a strict byte and advances one byte in the info. On the off chance that the hash chain is most certainly not

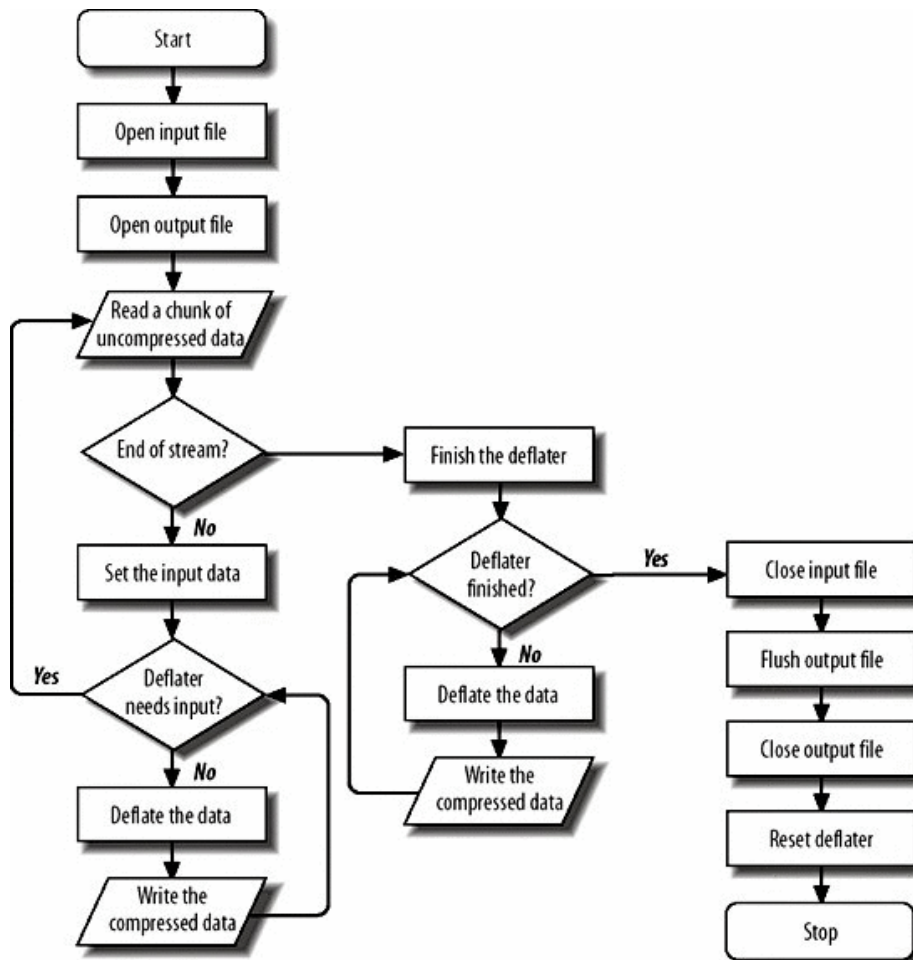


Fig 3-3: Flow chart of Deflate Algorithm

purge, showing that the succession XYZ (or, on the off chance that we are unfortunate, some other 3 bytes with a similar hash work esteem) has happened as of late, the compressor all strings on the XYZ hash chain with the genuine info information arrangement beginning at the present point, and chooses the longest match.

To enhance general compression, the compressor alternatively concedes the determination of matches ("sluggish coordinating"): after a match of length N has been found, the compressor looks for a more extended match beginning at the following information byte. On the off chance that it finds a more drawn out match, it truncates the past match to a length of one

(accordingly creating a solitary strict byte) and after that radiates the more extended match. Else, it radiates the first match, and, as depicted above, advances N bytes before proceeding. Run-time parameters additionally control this "lethargic match" technique. In the event that compression proportion is most imperative, the compressor endeavors an entire second hunt paying little mind to the length of the primary match. In the typical case, if the ebb and flow match is "sufficiently long", the compressor diminishes the scan for a more drawn out match, along these lines accelerating the procedure. In the event that speed is most imperative, the compressor embeds new strings in the hash table just when no match was found, or when the match is not "too long". This debases the compression proportion yet spares time since there are both less additions and less pursuits.

### 3.1 Prediction by partial matching

#### 3.1.1 Compression algorithm

- Calculate the conditional probability of each symbol using context models.
- Calculate its cumulative probability  $P_c$ .
- Begin with a current interval  $[L, H)$  initialized to  $[0, 1)$ .
- For each event in the file, perform two steps:
  - Subdivide the current interval into subintervals, one for each possible event. The size of an event's subinterval is proportional to the estimated probability that the event will be the next in the file, according to the model of the input.
  - Select the subinterval corresponding to the event that actually occurs next, making it the new current interval.
- Lower interval  $L = L + P_c * (H - L) * L$ ;
- Higher interval  $H = L + P_c * (H - L) * H$ ;
- Assign a unique identical tag for the message.

#### 3.1.2 Decompression algorithm

The tag value is taken; the new tag value is calculated using formula:

- $Tag = (Tag - L(s)) / P_c(s)$
- The corresponding lower interval and cumulative probability of symbol  $s$  is taken to decode the tag value. The process continues till the end of the text file. As the length of the source sequence increases, the length of the subinterval specified by the sequence decreases, and more bits are required to precisely identify the subinterval.

### 3.1.3 Model Development

PPM is a limited setting measurable displaying procedure which consolidates a few settled request setting models to foresee the following character in the information grouping. The forecast probabilities for every setting are adaptively refreshed from recurrence numbers. The most extreme setting length is a settled steady and has been found that expanding the length past 6 by and large does not enhance compression. The fundamental commence of PPM is to utilize the past bytes in the information stream to anticipate the accompanying one. Models that make their forecasts on a few promptly going before images are limited setting model of request  $k$ , where  $k$  is the measure of going before images utilized. PPM utilizes a few settled request setting models with various qualities beginning at 0 to some most extreme esteem. From each model, a different likelihood appropriation is gotten which are adequately consolidated into a solitary one where number juggling coding is utilized to encode the genuine character in respect to that dispersion. Escape probabilities are utilized for this mix to such an extent that if a setting can't be utilized for encoding an esteem, an escape image is transmitted and the model with the following littler estimation of  $k$  is utilized.

### 3.1.4 Analysis of prediction by partial matching

One of the principle qualities of the PPM calculation is that it is versatile. It gathers insights for a wide range of settings which makes the calculations memory escalated. The memory prerequisite of a clear PPM calculation execution is  $O(Mk+1)$  in the most pessimistic scenario, where  $M$  is the cardinality of the letters in order of information images, and  $K$  is the greatest setting length. Regarding the length of an information arrangement  $n$ , PPM calculations in the most pessimistic scenario require  $O(n^2)$  memory. After accepting and handling each image from the information stream, the PPM model is refreshed. On the off chance that the up and coming information image exists in the current PPM show, its recurrence is refreshed. For the situation when the up and coming image is new to the present setting, another passage is made in every pertinent setting. The refresh of the PPM model is non-specific. Notwithstanding the recurrence of the further utilization of the made sections, they influence the likelihood evaluations of whatever remains of the forthcoming images. Due to the down to earth confinements of Arithmetic and Logic Units (ALUs), image frequencies are rescaled when the biggest recurrence achieves a specific farthest point. This enhances the territory of the PPM show on a present segment of the info information stream. Rescaling likewise decreases the impact of the from time to time happening images on the general likelihood appraise.

## Chapter-4

### PERFORMANCE ANALYSIS

Execution analysis of compression calculations should be possible by different variables. Be that as it may, the principle concern has dependably been the space effectiveness and time proficiency. We are utilizing diverse variables to break down the calculation.

#### 4.1 Compression Ratio

Compression proportion, otherwise called compression power, is utilized to evaluate the lessening in information portrayal estimate created by an information compression calculation. The information compression proportion is undifferentiated from the physical compression proportion used to gauge physical compression of substances.

Information compression proportion is characterized as the proportion between the uncompressed estimate and compacted measure.

$$\text{Compression Ratio} = \frac{\text{Uncompressed Size}}{\text{Compressed Size}}$$

In this manner a portrayal that packs a 10 MB record to 2 MB has a compression proportion of  $10/2 = 5$ , frequently documented as an express proportion, 5:1 (read "five" to "one"), or as an understood proportion, 5/1. Take note of that this detailing applies similarly for compression, where the uncompressed size is that of the first; and for decompression, where the uncompressed size is that of the proliferation.

Some of the time the space reserve funds is given rather, which is characterized as the decrease in size with respect to the uncompressed estimate:

$$\text{Space Savings} = 1 - \frac{\text{Compressed Size}}{\text{Uncompressed Size}}$$

In this way a portrayal that packs a 10MB record to 2MB would yield a space reserve funds of  $1 - 2/10 = 0.8$ , regularly documented as a rate, 80%.

For signs of uncertain size, for example, spilling sound and video, the compression proportion is characterized as far as uncompressed and packed information rates rather than data sizes:

$$\text{Compression Ratio} = \frac{\text{Uncompressed Data Rate}}{\text{Compressed Data Rate}}$$

Despite of space savings, one speaks of data-rate savings, which is defined as

$$\text{Data Rate Savings} = 1 - \frac{\text{Compressed Data Rate}}{\text{Uncompressed Data Rate}}$$

For instance, uncompressed tunes in CD arrange have an information rate of 16 bits/channel x 2 channels x 44.1 kHz  $\cong$  1.4 Mbit/s, though AAC records on an iPod are regularly packed to 128 Kbit/s, yielding a compression proportion of 10.9, for an information rate investment funds of 0.91, or 91%.

At the point when the uncompressed information rate is known, the compression proportion can be deduced from the compacted information rate.

#### 4.2 Compression Speed

Compression speed is identified with the information design and the machine sort. The connection between application execution and host machine parameters is an exploration theme that is outside of the extent of this paper. Amid the tests, we continue utilizing a similar machine for every one of the compressions, and ensure that our application is the main workload. Along these lines, we can consider compression speed as an element of compression calculation. The compression speed is likewise influenced by compression cradle measure, however we preclude this component by utilizing a similar size of support, which is 16KB.

While assessing information compression calculations, speed is dependably as far as uncompressed information dealt with every second.

A few applications utilize information compression methods notwithstanding when they have so much RAM and plate space that there's no genuine need to make documents littler. Document compression and delta compression are regularly used to accelerate duplicating records from one end of an ease back association with another. Indeed, even on a solitary PC, a few sorts of operations are altogether speedier when performed on packed adaptations of information as opposed to straightforwardly on the uncompressed information. Specifically, some compacted document organizations are outlined so that packed example coordinating - hunting down an expression in a compacted rendition of a content record - is altogether quicker than scanning for that same expression in the first uncompressed content document.

$$\text{speed} = \frac{\text{Uncompressed bits}}{\text{seconds to compress}}$$

In a couple of utilizations, the compression speed is basic. On the off chance that a specific usage of a sound compressor running on a model voice recorder can't maintain 7 bits/test/channel x 1 channel x 8 kSamples/s = 56 kbit/s from the mouthpieces to capacity, then it is unusable. Nobody needs their recorded voice to have noiseless holes where the equipment couldn't keep up. Nobody will get it unless you change to an alternate execution or quicker equipment (or both) that can stay aware of standard phone quality voice speeds.

The speed changes generally starting with one machine then onto the next, starting with one usage then onto the next. Indeed, even on a similar machine and same benchmark document and same usage source code, utilizing an alternate compiler may make a decompressor run faster. The speed of a compressor is quite often slower than the speed of its comparing decompressor.

Indeed, even with a quick current CPU, packed document framework execution is regularly restricted by the speed of the compression calculation. Numerous advanced installed frameworks - and additionally a hefty portion of the early PCs that information compression calculations were initially created on - are intensely compelled by speed.

#### 4.3 Compression Time

The time taken by the algorithm to compress the file. Calculated in milliseconds (ms).

<b>Test Data</b>	<b>File Size (bytes)</b>	<b>Compressed File Size(bytes)</b>	<b>Space Saving (%)</b>	<b>Compression Time (min.)</b>	<b>Compression Ratio</b>
<b>Sources</b>	52428800	42656072	18.64	200.3	1.22
<b>English</b>	52428800	38666240	26.25	160.32	1.35
<b>Pitches</b>	52428800	37722521	28.05	150.89	1.38
<b>Proteins</b>	52428800	34807480	33.61	180.15	1.50
<b>Dna</b>	52428800	39075184	25.47	190.30	1.34
<b>Xml</b>	52428800	46299873	11.69	260.62	1.13

Table 2: Results of Arithmetic Coding



Test Data	File Size (bytes)	Compressed File Size(bytes)	Space Saving (%)	Compression Time (min.)	Compression Ratio
Sources	52428800	36329886	30.70	191.8	1.44
English	52428800	27895124	46.79	95.7	1.87
Pitches	52428800	34173092	34.82	188.1	1.53
Proteins	52428800	32691732	37.64	174	1.60
Dna	52428800	14122728	73.06	84.4	3.71
Xml	52428800	41230008	21.36	245.3	1.27

Table 3: Results of LZW algorithm

Test Data	File Size (bytes)	Compressed File Size(bytes)	Space Saving (%)	Compression Time (sec.)	Compression Ratio
Sources	52428800	95180898	-	64.136	0.55
English	52428800	101898190	-	73.039	0.51
Pitches	52428800	89272096	-	127.150	0.58
Proteins	52428800	97175532	-	76.864	0.53
Dna	52428800	73564246	-	56.975	0.71
Xml	52428800	102465536	-	71.616	0.51

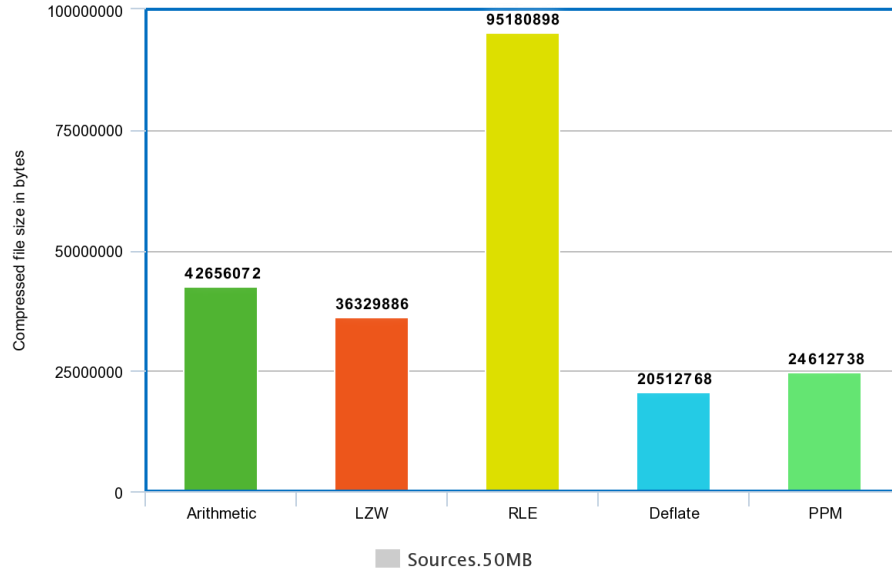
Table 4: Results of Run Length Encoding

Test Data	File Size (bytes)	Compressed File Size(bytes)	Space Saving (%)	Compression Time (sec.)	Compression Ratio
Sources	52428800	20512768	60.87	66.214	2.55
English	52428800	36405248	30.56	88.093	1.44
Pitches	52428800	23363584	55.43	77.528	2.24
Proteins	52428800	41353216	21.12	133.639	1.26
Dna	52428800	31162368	40.56	80.032	1.68
Xml	52428800	21312768	59.34	51.295	2.45

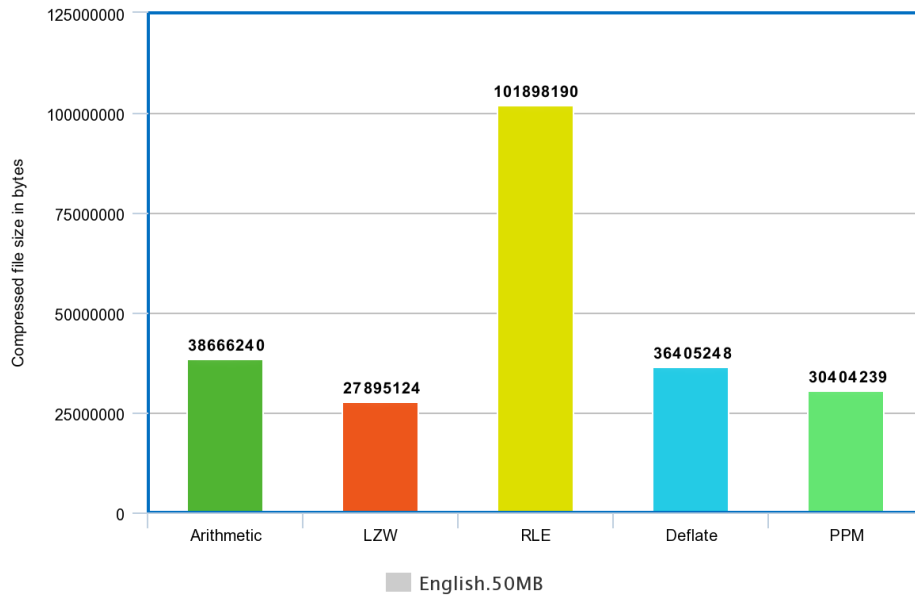
Table 5: Results of Deflate algorithm

Test Data	File Size (bytes)	Compressed File Size(bytes)	Space Saving (%)	Compression Time (min.)	Compression Ratio
Sources	52428800	24612738	53.05	221.3	2.13
English	52428800	30404239	42.01	150.35	1.72
Pitches	52428800	22353674	57.36	120.46	2.34
Proteins	52428800	35403516	32.47	151.33	1.48
Dna	52428800	29142345	44.41	210.12	1.79
Xml	52428800	20992738	59.95	250.41	2.49

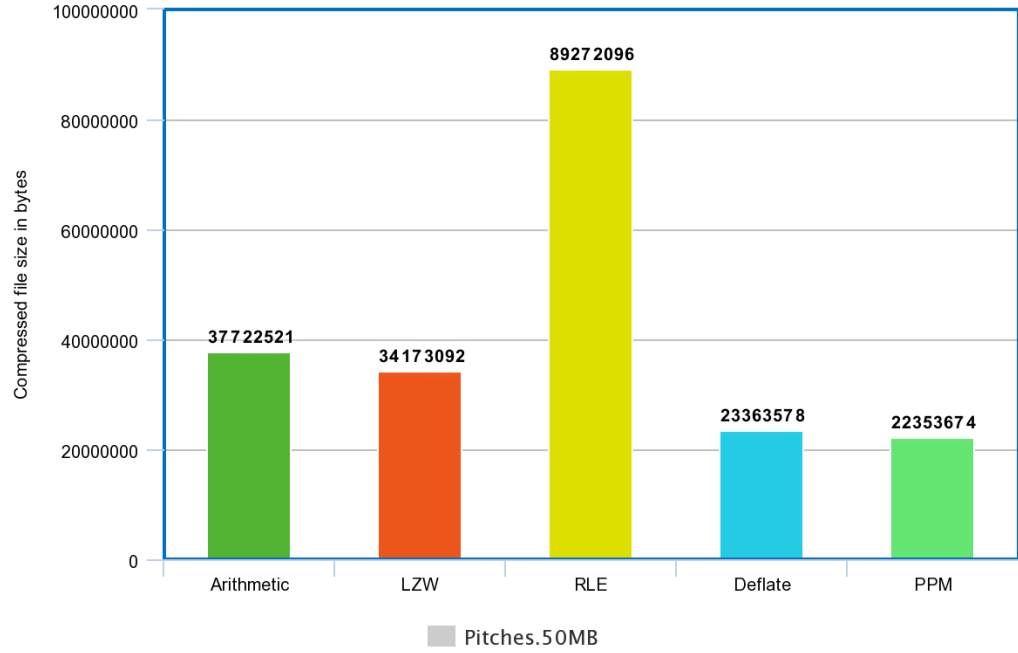
Table 6: Results of Prediction by partial matching algorithm (PPM)



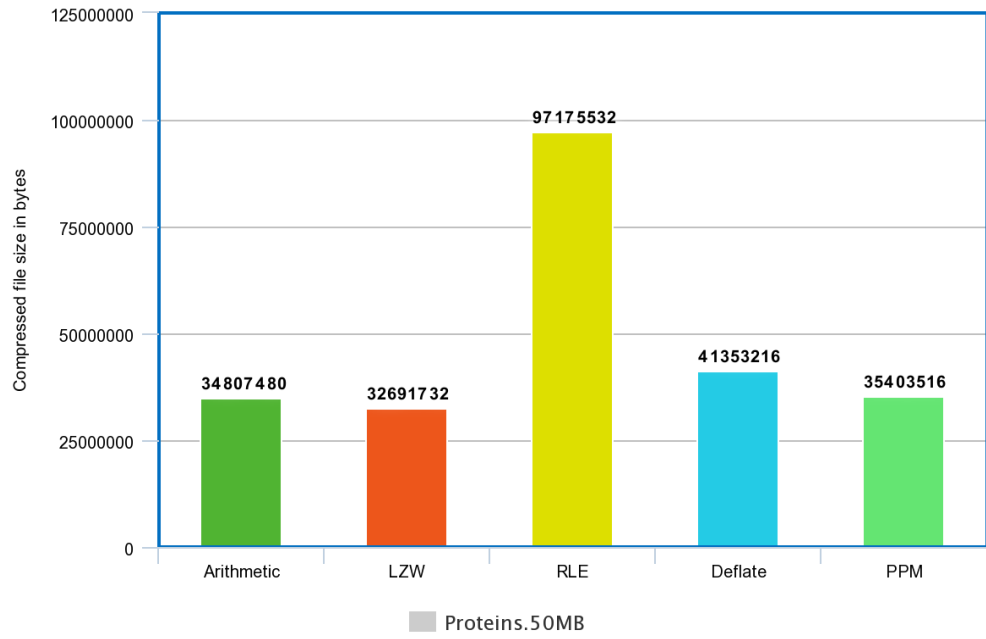
Graph 1: Comparison of compressed file size of Sources.50MB



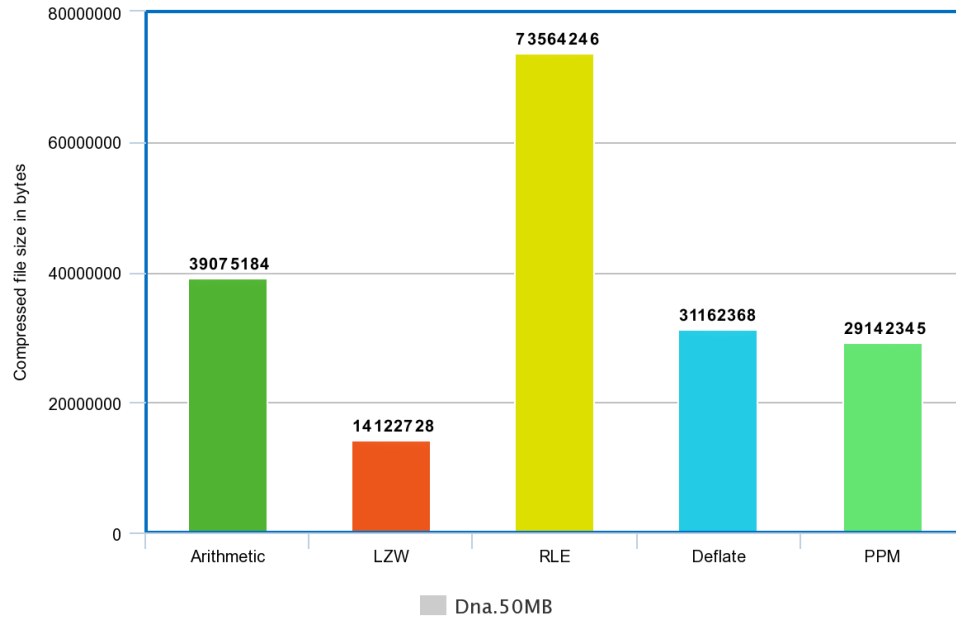
Graph 2: Comparison of compressed file size of English.50MB



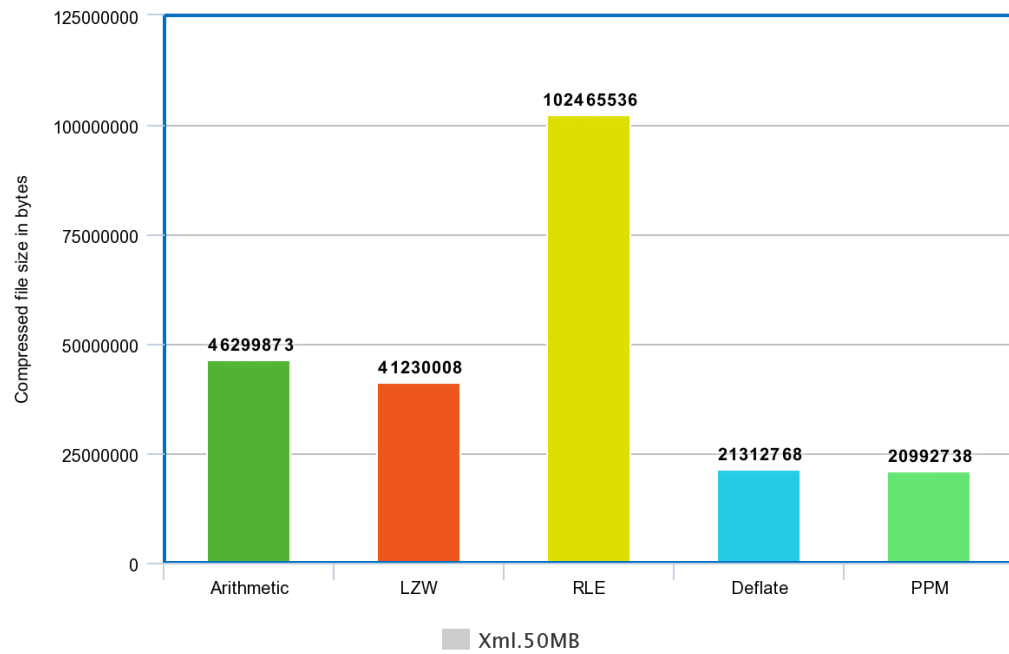
Graph 3: Comparison of compressed file size of Pitches.50MB



Graph 4: Comparison of compressed file size of Proteins.50MB



Graph 5: Comparison of compressed file size of Dna.50MB



Graph 6: Comparison of compressed file size of Xml.50MB

## Chapter-5

### CONCLUSION

#### 5.1 Conclusion

We have demonstrated the subtle elements of an execution of number-crunching coding and have called attention to its focal points (adaptability and close optimality) and its primary detriment (slowness). We have built up a quick coder, in view of decreased exactness math coding, which gives just insignificant loss of compression effectiveness; we can utilize the idea of  $\epsilon$ -parcels and the probabilities to incorporate into the coder to keep the compression misfortune little.

Underneath we list the primary conclusions coming about because of our analyses.

- There is a generous speed pick up as we move from renormalizations that spare one piece a period, to those that spare bits together in gatherings of at least one bytes.
- Byte-based renormalizations require enough accuracy from the number juggling operations (e.g., no less than 16 or 32 bits) to bolster a more extensive scope of interim lengths. Typically these can be best bolstered by the local CPU operations, rather than approximations.
- Multiplications are currently adequately quick, and their effect on the coding rate is little notwithstanding for static double coders.
- While double coders play out all the coding operations in the briefest time, their data throughput is restricted to at most one piece for each coded image. For quickest coding we ought to utilize strategies that code images from bigger letter sets since they can yield considerably higher throughputs.
- Arithmetic disentangling can be altogether slower than encoding, due to the look for the interim to which the coded image has a place. The best arrangement relies on upon the processor and information source.

This calculation packs tedious groupings of information well.

Since the code words are 12 bits, any single encoded character will extend the information estimate as opposed to lessen it.

72 bits are spoken to with 72 bits of information. After a sensible string table is constructed, compression enhances drastically.

Preferences of LZW over Huffman:

- LZW requires no earlier data about the info information stream.
- LZW can pack the info stream in one single pass.
- Another favorable position of LZW its effortlessness, permitting quick execution.

Restrictions of LZW:

- What happens when the word reference gets too extensive (i.e., when all the 4096 areas have been utilized)?
- Here are a few alternatives normally actualized:
  - o Simply disregard including any more sections and utilize the table as seems to be.
  - o Throw the word reference away when it achieves a specific size.
  - o Throw the word reference away when it is no longer compelling at compression.
  - o Clear passages 256-4095 and begin fabricating the word reference once more.
- Some cunning plans remake a string table from the keep going N input character.

## 5.2 Future Scope

On the off chance that the source entropy is little (e.g., underneath 3 bits/image), then it is most likely best to utilize a hunt strategy that utilizations just augmentations, and attempt to advance the inquiry succession [1, 2].

Despite the fact that division is slower than alternate operations, we can maintain a strategic distance from long hunts by utilizing one division for each decoded image and a table turn upward to initialize the inquiry. Little tables can altogether accelerate the hunt.

As the future work on compression of pictures for putting away and transmitting should be possible by different lossless techniques for picture compression in light of the fact that as it is finished up over, that the consequence of the decompressed picture is practically same as that of the information picture so it demonstrates that there is no loss of data amid transmission. So different strategies for picture compression, any of the sort i.e lossless or lossy can be done as to be specific JPEG strategy and so forth. Utilization of various measurements can likewise occur to assess the execution of compression calculations.

Computerized media substance can be made, altered, appropriated, imparted and put away to accommodation effortlessly. Since the rising wired and remote IP systems are open systems, they are helpless against spying and classification is particularly critical for secure media circulations over IP systems. Ordinary cryptographic methods created for ensuring content information can't be connected for encoding the whole media stream because of assortment of requirements, for example, the computational overhead, computational cost, level of multifaceted nature and so on., which is not satisfactory for ongoing applications. Then again, utilizing quick yet shaky scrambling strategies is likewise not worthy by and large. In this manner, the ebb and flow research is centered around changing and upgrading the current cryptosystems appropriate for ongoing mixed media.

Promising future headings of research incorporate more accentuation on the accompanying zones:

- Key administration is a basic issue in all encryption based security frameworks, as it can't be isolated from the outline of secure media circulation. In most conveyance models, mixed media substance is scrambled with a symmetric key which additionally should be ensured in transmission to the recipient. Henceforth, the capacity and security prerequisites of key administration should be examined in more noteworthy detail in future recommendations.
- Another include that might be added to the proposed particular plans is the choice criteria. Encryption systems can be picked powerfully as the substance is being dispersed and the choice criteria can be changed as required by the application.
- Enhancement in compression execution by presentation of new functionalities which additionally enhances security as encryption is joined with compression.

### 5.3 Applications Contributions

It is utilized generally for often happening arrangements of pixels. Math coding can be utilized for application to information compression for VLSI testing. The utilization of math codes brings about a code word whose length is near the ideal esteem (as anticipated by entropy in data hypothesis), accordingly accomplishing a higher compression. Past strategies, (for example, those in view of Huffman or Golomb coding) result in ideal codes for informational indexes in which the likelihood model of the images fulfills particular necessities. This paper indicates observationally and logically that Huffman and Golomb codes can bring about a vast contrast between the bound set up by the entropy and the achieved compression; hence, the most pessimistic scenario distinction is contemplated utilizing data hypothesis. Compression comes about for number juggling coding are introduced utilizing ISCAS benchmark circuits; a viable whole number usage of math coding/deciphering and an investigation of its deviation from the entropy bound are sought after. A product execution is proposed utilizing installed DSP centers. In the test assessment, completely determined test vectors and test 3D shapes



from two distinctive ATPG projects are used. The ramifications of math coding on assembling test utilizing an ATE are likewise explored.

With the headways in compression innovation, it is currently simple and effective to pack video documents. Different video compression procedures are accessible. The most well-known video compression standard is MPEG (Moving Picture Experts Group). It is a working gathering of ISO/IEC accused of the advancement of video and sound encoding principles. MPEG's authentic assignment is ISO/IEC JTC1/SC29 WG11. Many advances are being made for enhancing the video quality. Advancements in MPEG standard are MPEG-1 (MP3), MPEG-2, MPEG-3, MPEG-4 (MPEG-4 Part 2 or Advanced Simple Profile) and MPEG-4 Part 10 (or Advanced Video Coding or H.264). MPEG-7. A formal framework for depicting media content. MPEG-21 portrays this standard as a sight and sound structure. MPEG standard is extremely effective in the utilization of DVDs. Different H.261 benchmarks could be utilized as a part of future for headway in video conferencing innovation. Other connected fields that are making utilization of wavelets in the coming future, incorporate stargazing, acoustics, atomic designing, sub-band coding, flag and picture handling, neurophysiology, music, attractive reverberation imaging, discourse segregation, optics, fractals, turbulence, seismic tremor expectation, radar, human vision, and immaculate arithmetic applications, for example, understanding incomplete differential conditions.

They can likewise be utilized as a part of TIFF, GIF and PDF documents. In late applications LZW has been supplanted by the more effective Flate calculation.

The Proposed Algorithm keeps away from a considerable lot of the issues related with different techniques for Compression in that it powerfully adjusts to the Redundancy characteristics of the information being packed. The Effectiveness of Compression is communicated as a proportion relating character in the quantity of bits expected to express the message prior and then afterward Compression. R. Nigel Horspool, the creator portrays a basic approach to enhance the compression without altogether corrupting its speed is proposed, and exploratory information demonstrates that it works by and by and even better outcomes are accomplished. The Lempel-Ziv-Welch (LZW) compression calculation is broadly utilized in light of the fact that it accomplishes an incredible trade off between compression execution and speed of execution. Check Daniel Ward, the Author researches the LempelZiv '77 information compression calculation by considering calculation for productively implanting strings in paired trees and Analysis of the assortment coordinating parameter of addition trees was likewise introduced.

## REFERENCES

- [1] Khalid Sayood, "Introduction to Data Compression", Ed Fox (Editor), March 2000.
- [2] Burrows M., and Wheeler, D. J. 1994," A Block-Sorting Lossless Data Compression Algorithm" SRC Research Report 124, Digital Systems Research Center.
- [3] C.E. Shannon, "A mathematical theory of communication," Bell Syst. Tech. J., vol. 27, pp. 398-403.
- [4] Glen G. Langdon, Jr, "An Introduction to Arithmetic Coding." IBM Research Division, California.
- [5]"Data Compression Methodologies for LossLess Data and Comparison between Algorithms", IJESIT Volume 2, Issue 2, March 2013.
- [6] Amir Said, "Introduction to Arithmetic Coding - Theory and Practice", Imaging Systems Laboratory, 2004.
- [7] Somefun, M. Adebayo & Adewale, "Evaluation of dominant text data compression techniques," IJAIEM, 2014.
- [8] I.H. Witten, R.M. Neal, and J.G. Cleary, "Arithmetic Coding for the data compression," Commun. ACM, vol. 30, no. 6, pp. 520-540, June 1987.
- [9] R. Pasco," Source coding algorithms for fast data compression," Stanford Univ., Ph.D. dissertation, 1976.
- [10] J.J. Rissanen, "Generalized Kraft inequality and arithmetic coding," IBM J. Res. Devel. , vol. 20, no. 3, pp. 198-203, May 1976.
- [11] F. Rubin, " Arithmetic stream coding using fixed precision registers," IEEE Trans. Information Theory, vol. IT-25, no. 6, pp. 520-540, June 1987.
- [12] J.J. Rissanen and G.G. Langdon , " Arithmetic coding," IBM J. Res. Devel, vol. 23 no. 2, pp. 146-162, Mar. 1979.
- [13] M. Guazoo, "A general minimum-redundancy source-coding algorithm," IEEE Trans. Information Theory, vol. IT-26, no. 1, pp. 15-25, Jan 1980.
- [14] Manjeet Kaur, Er. Upasna Garg," Lossless Text Data Compression Algorithm Using Modified Huffman Algorithm," IJARCSSE, vol. 5, Issue. 7, 2015.
- [15] A. Said, "Comparative Analysis of Arithmetic Coding Computational Complexity," Hewlett Packard Laboratories Report, HPL-2004-75, Palo Alto, CA, April 2004.

- [16] M. Schindler, "A fast renormalization for arithmetic coding," Proc. IEEE Data Compression Conf., 1998.
- [17] Texas Instruments Incorporated, "TMS320C6000 CPU and Instruction Set Reference Guide," Literature Number: SPRU189F, Dallas, TX, 2000.
- [18] International Business Machines Corporation, "PowerPC 750CX/CXe RISC Microprocessor User's Manual," (preliminary edition), Hopewell Junction, NY, 2001.
- [19] Intel Corporation, "Intel Pentium 4 Processor Optimization," Reference Manual 248966, Santa Clara, CA, 2001.
- [20] Sun Microsystems Inc., "UltraSPARC III Technical Highlights," Palo Alto, CA, 2001.
- [21] Welch, Terry (1984), "A Technique for High-Performance Data Compression". (6): 8–19. doi:10.1109/MC.1984.1659158.
- [22] Jump up ^ Ziv, J.; Lempel, A. (1978). "Compression of individual sequences via variable-rate coding", IEEE Transactions on Information Theory. 24 (5): 530. doi:10.1109/TIT.1978.1055934
- [23] D.S. Taubman and M.W. Marcellin, "JPEG 2000: Image Compression Fundamentals," Standards and Practice, Kluwer Academic Publishers, Boston, MA, 2002.
- [24] P. Deutsch, "DEFLATE Compressed Data Format", Aladdin Enterprises Category: Informational May 1996.
- [25] R. Pasco, "Source coding algorithms for fast data compression," Ph.D. dissertation, Stanford University, 1976.
- [26] K. Sayood, "Lossless compression handbook," Academic Press, 2003.
- [27] P. Deutsch., "Deflate compressed data format specification version 1.3," RFC 1951, <http://www.faqs.org/rfcs/rfc1951.htm>, 1996.
- [28] Hiroyuki, A., Kazuhiro, K., Takashi, I., Shigeichi, H., "A PPM\* algorithm using context mixture," In the Journal of IEIC, 2005.
- [29] Moffat, A., "Implementing the PPM data compression scheme," In IEEE Transactions on Communications, 1990.
- [30] Teahan, W. J., Harper, D. J., "Combining PPM models using a text mining approach," In proceedings of IEEE Data Compression Conference, 2001.

[31] Gupta G., Gupta K.L. Jyoti A. , “AN ADVANCED COMPRESSION APPROACH WITH RLE FOR IMAGE COMPRESSION” International Journal of Advanced Research in Computer Science and Software Engineering Volume 4, Issue 2, February 2014.

[32] Nagarajan A., Alagarsamy K. “AN ENHANCED APPROACH IN RUN LENGTH ENCODING SCHEME” International Journal of Engineering Trends and Technology- July to Aug Issue 2011.

[33] Amin A., Aheman Q. ,Junaid M. Habib M.Y, Anjum W. , “MODIFIED RUN LENGTH ENCODING SCHEME WITH INTRODUCTION OF BIT STUFFING FOR EFFICIENT DATA COMPRESSION” 6th International conference on internet technology and secured transaction 11-14 December 2011 Abu Dhabi (978-1-908320-00-1- /11/\$26.00 @ 2011 IEEE).

[34] Akhtar M.B., Qureshi A.M. , and Islam Q., “OPTIMIZED RUN LENGTH CODING FOR JPGE IMAGE COMPRESSION USED IN SPEC RESEARCH PROGRAM OF IST.( 978-1-61284-941-6/11/\$26.00 ©2011 IEEE).

[35] Joseph S., Srikanth N. “A NOVEL APPROACH OF MODIFIED RUN LENGTH ENCODING SCHEME FOR HIGH SPEED DATA COMMUNICATION APPLICATION” International Journal of Science and Research (IJSR) Volume 2 Issue 12, December 2013.

#### Web References:

[http://www.stringology.org/DataCompression/ak-int/index\\_en.html](http://www.stringology.org/DataCompression/ak-int/index_en.html)

[http://akbar.marlboro.edu/~mahoney/courses/Fall01/computation/compression/ac/ac\\_arithmetic.html](http://akbar.marlboro.edu/~mahoney/courses/Fall01/computation/compression/ac/ac_arithmetic.html)

<http://cotty.16x16.com/compress/nelson1.htm>

<http://www.drdoobs.com/parallel/arithmetic-coding-and-statistical-model/184408491>

<http://www.dspguide.com/ch27/5.htm>

<https://www.cs.duke.edu/csed/curious/compression/lzw.html>

<http://pizzachili.dcc.uchile.cl/texts.html>