

# FINAL PROJECT REPORT v1.0

*A Statistical analysis and extended implementation of  
Compression algorithms*

Guided By: Prof. Bhaskar Chaudhary

Prepared By: Rachit Mishra (201201092)

Dhrumil Shah (201201117)

## *Abstract*

*This project report presents a detailed analysis of some compression algorithms and their implementation with some new functionalities developed for them in an attempt to increase their efficiency. While conducting the research for this project, a systematic procedure of testing was done on these compression algorithms and the details on a comparative analysis of results and properties exhibited by different algorithms under different parameters has been tabulated in this report in both statistical and graphical forms.*

## INTRODUCTION

- ❖ Data Compression as the name suggests is a technique in by which the same amount of **data** is transmitted by using a smaller number of bits; for example, by replacing a string of twenty repeated digits with a command to repeat the digit twenty times.
- ❖ Data Compression can be either lossy or lossless. As the names suggest, in lossless data compression the information is not lost and all of the information at the beginning is as it is at the end. It reduces bits by identifying and eliminating statistical redundancy.
- ❖ Data Compression shrinks down a file so that it takes up less space. This is desirable for data storage and data communication.
- ❖ Storage space on disks is expensive so a file which occupies less disk space is "cheaper" than an uncompressed file. Smaller files are also desirable for data communication, because the smaller a file the faster it can be transferred.
- ❖ A compressed file appears to increase the speed of data transfer over an uncompressed file.
- ❖ Data Compression shrinks down a file so that it takes up less space. This is desirable for data storage and data communication.
- ❖ Storage space on disks is expensive so a file which occupies less disk space is "cheaper" than an uncompressed file. Smaller files are also desirable for data communication, because the smaller a file the faster it can be transferred.
- ❖ A compressed file appears to increase the speed of data transfer over an uncompressed file.

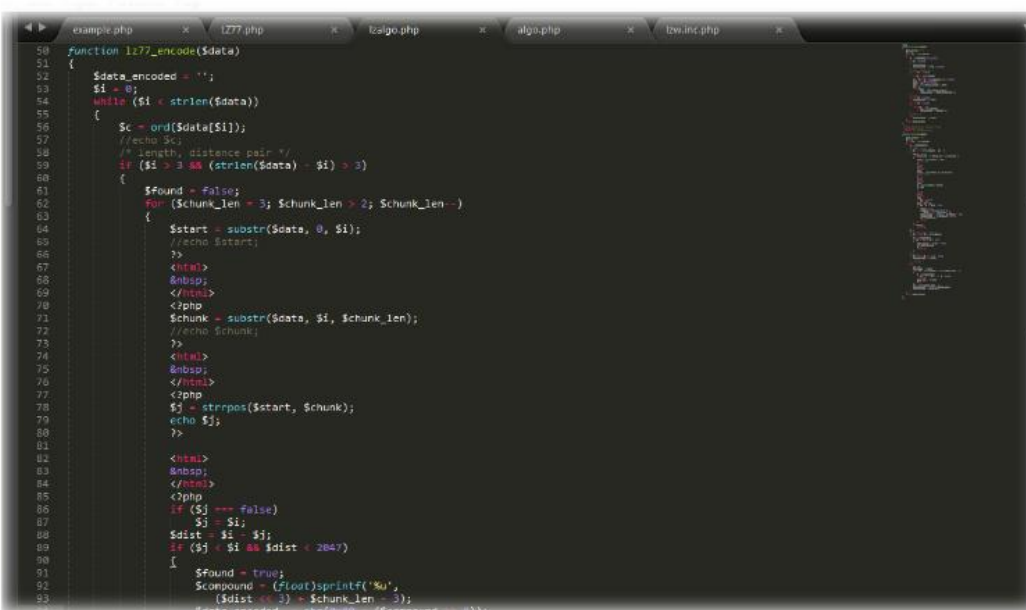
## ALGORITHMS COVERED

The algorithms we worked upon and which we have covered in these report are:

### 1. Huffman encoding

## 2. LZW algorithm (Precisely, we implemented LZ 77 which is an extended variation of the LZW algorithm.

- LZ77 and LZ78 are lossless data compression algorithms : Forms basis of **DEFLATE** used in PNG files) and GIF
- Lossless Data Compression algorithms: Original data is completely reconstructed from the compressed data. Lossy data compression allows compression of only an approximation of the original data.(although improves compression rates)
- LDC algorithms used in ZIP file formats. PNG and GIF use lossless compressions. 3D graphics (OpenCTM), Cryptography.
- Both LZ77 and LZ78 are **DICTIONARY CODERS.(Substitution coders)** : A class of LDC algorithms.
- Look for matches between data set to be compressed and a set of strings (contained in a data structure maintained by the encoder) called dictionary.
- When a match is found, a reference to the string's position in the data structure is given.
- In LZ77 compression MUST start at the beginning of the input itself.
- Typically every character is stored with 8 bits (up to 256 unique symbols for this data)..... Can be extended to 9 to 12 bits per character. New unique symbols made of combinations of symbols which occurred previously in the string.



```

50 function lz77_encode($data)
51 {
52     $data_encoded = '';
53     $i = 0;
54     while ($i < strlen($data))
55     {
56         $c = ord($data[$i]);
57         //echo $c;
58         /* length, distance pair */
59         if ($i > 3 && (strlen($data) - $i) > 3)
60         {
61             $found = false;
62             for ($chunk_len = 3; $chunk_len > 2; $chunk_len--)
63             {
64                 $start = substr($data, 0, $i);
65                 //echo $start;
66                 <?php
67                 <html>
68                 &nbsp;
69                 </html>
70                 <?php
71                 $chunk = substr($data, $i, $chunk_len);
72                 //echo $chunk;
73                 <?php
74                 <html>
75                 &nbsp;
76                 </html>
77                 <?php
78                 $j = strpos($start, $chunk);
79                 echo $j;
80                 <?php
81                 <html>
82                 &nbsp;
83                 </html>
84                 <?php
85                 if ($j === false)
86                 {
87                     $i = $i;
88                     $dist = $i - $j;
89                     if ($j < $i && $dist < 2047)
90                     {
91                         $found = true;
92                         $compound = (float)sprintf('%u',
93                             ($dist < 3) ? $chunk_len : 3);
94                         $data_encoded .= chr(0x00) . ($compound < 3);
95                     }
96                 }
97             }
98         }
99     }
100 }

```

Below 2 screenshots depicts the code for compressing strings via LZ77 algorithm.

```

92     $compound = (float)sprintf('%u',
93         ($dist << 3) + $chunk_len - 3);
94     $data_encoded .= chr(0x00 + ($compound >> 8));
95     $data_encoded .= chr($compound & 0xFF);
96     $i += $chunk_len;
97     break;
98 }
99 }
100 if ($found)
101     continue;
102 }
103 $i++;
104 /* byte pair (space + char) */
105 if ($c == 32 && $i < strlen($data))
106 {
107     $n = ord($data[$i]);
108     if ($n >= 0x40 && $n <= 0x7F)
109     {
110         $data_encoded .= chr($n ^ 0x80);
111         echo $data_encoded;
112         $i++;
113         continue;
114     }
115 }
116 /* literal */
117 if ($c == 0 || ($c > 8 && $c < 0x80))
118     $data_encoded .= chr($c);
119 /* c literals */
120 else
121 {
122     $j = $i;
123     $data_temp = chr($c);
124     while ($j < strlen($data) && strlen($data_temp) < 8)
125     {
126         $c = ord($data[$j]);
127         if ($c == 0 || ($c > 8 && $c < 0x80))
128             break;
129         $data_temp .= chr($c);
130         $j++;
131     }
132     $i += strlen($data_temp) - 1;
133     $data_encoded .= chr(strlen($data_temp));
134     $data_encoded .= $data_temp;
135 }
136 }
137 }
138 return $data_encoded;
139 }

```

The one screenshot shown below describes the method for decompressing strings via LZ77 algorithm.

```

1 <?php
2 function lz77_decode($data)
3 {
4     $data_decoded = '';
5     $i = 0;
6     while ($i < strlen($data))
7     {
8         $c = ord($data[$i]) & 0xFF;
9         /* byte pair (space + char) */
10        if ($c >= 0x0000)
11        {
12            $data_decoded .= chr($c & 0xFF);
13        }
14        /* length, distance pair */
15        else if ($c >= 0x0080)
16        {
17            if ($i >= strlen($data))
18                continue;
19            $c = ($c << 8) | (ord($data[$i+1]) & 0xFF);
20            $len = ($c >> 8) & 0xFF;
21            $dist = ($c >> 3) & 0x7FF;
22            $pos = strlen($data_decoded) - $dist;
23            $j = 0;
24            while ($len-- > 0)
25            {
26                if ($pos < strlen($data_decoded))
27                    $data_decoded .= $data_decoded[$pos++];
28            }
29            /* literal */
30            else if ($c >= 0x0000)
31            {
32                $data_decoded .= chr($c);
33            }
34            /* c literals */
35            else if ($c >= 0x0001)
36            {
37                while ($c-- > 0)
38                {
39                    if ($i < strlen($data))
40                        $data_decoded .= $data[$i++];
41                }
42            }
43            else
44            {
45                $data_decoded .= chr($c);
46            }
47        }
48        $i++;
49    }
50    return $data_decoded;
51 }
52
53 /*
54  * Encode string to LZ77 compressed string.
55  * @param $data string.
56  * @return LZ77 compressed string.
57 */

```

### 3. Run Length Encoding algorithm

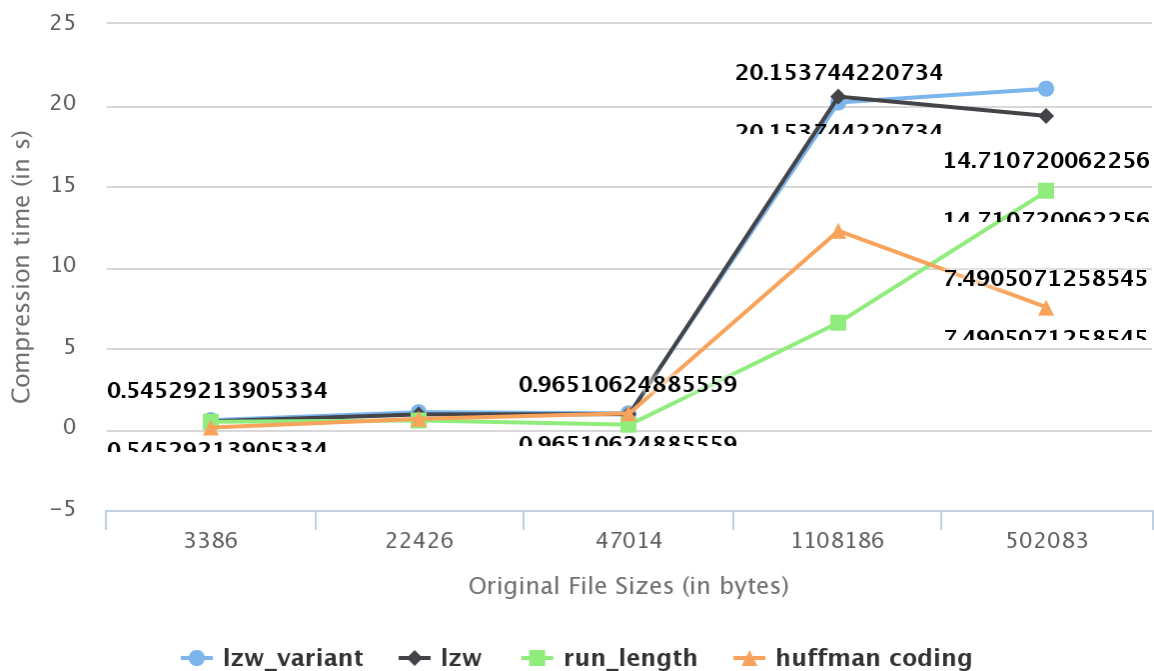
#### 4. Arithmetic Encoding algorithm

### GRAPHICAL ANALYSIS

Based on few parameters, and after implementing all the algorithms mentioned above, we generated graphs highlighting the significant comparisons between these algorithms based on some important factors like compression rates, time taken to compress the files and some other factors.

- Factor describing the original file size versus the size of the compressed file in case of all four algorithms.

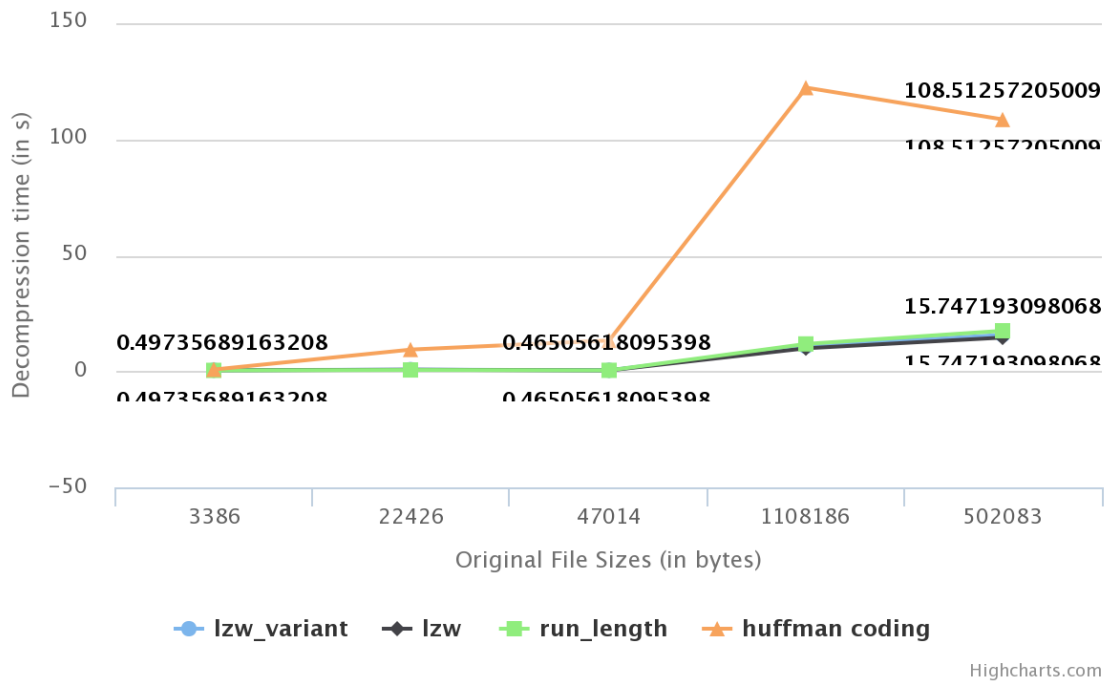
Compression time Vs original File



Highcharts.com

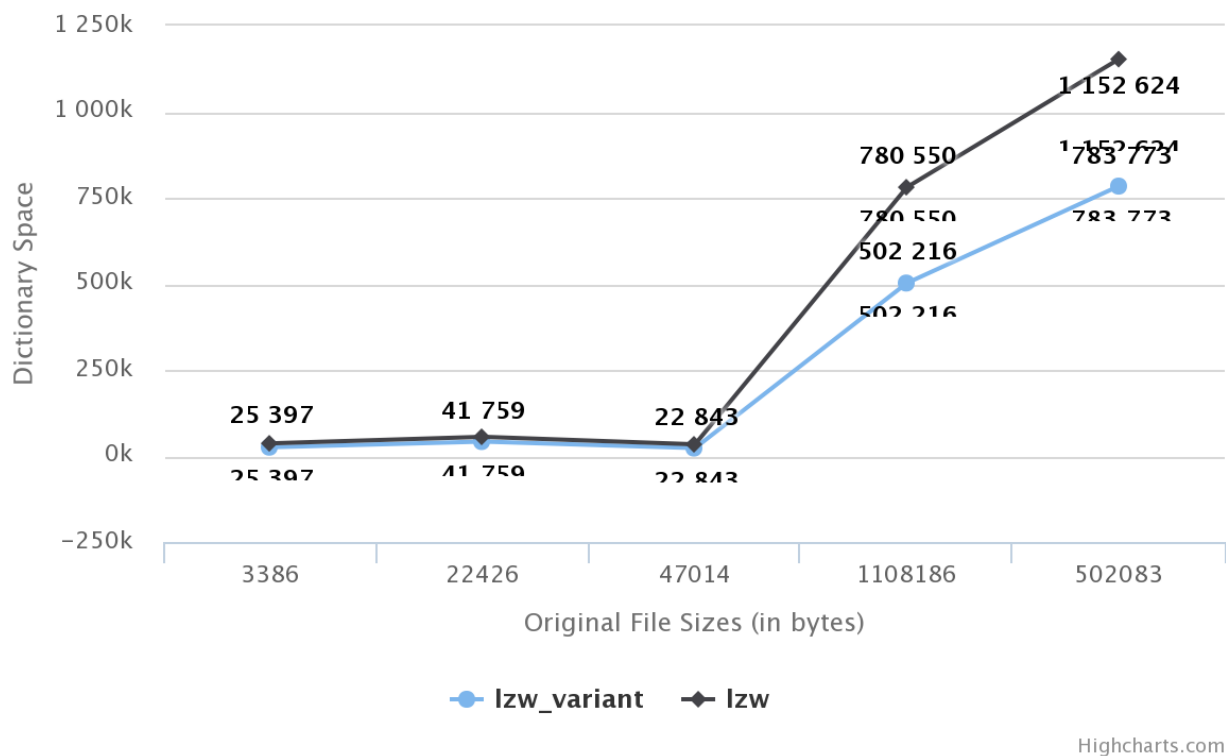
- Factor describing the decompression time versus original file

## Decompression time Vs original File



3. Dictionary space utilized for to compress the original file with some file size.

## Memory space Vs original File



#### 4. Compressed files vs Original file size comparison for different algorithms.

