

Lempel Zip Compression

Victor Rivera
Early Design Document

1 Lempel Zip Compression Design Goal

The goal of this project was to create two programs which perform Jacob Ziv's LZ78 compression and decompression algorithms. The program algorithms are lossless algorithms that allows users to compress and decompress text, binary files, and source code.

2 Design Overview

2.1 Program Files

This project contains five files used for encoding and decoding, several test scripts, a Makefile, this Design PDF, and a README.md.

2.1.1 Encoding/Decoding

The compressing and decompressing program files were all written in C and include the following:

- [trie.c/trie.h](#)
- [word.c/word.h](#)
- [io.c/io.h](#)
- [code.h](#)
- [endian.h](#)
- [encode.c](#)
- [decode.c](#)

2.1.2 Supporting Files

- [Makefile](#) (File used to create and delete encode and decode binary executables)
- [test files](#) (Several test scripts written in Bash)

2.1.3 Design Files

- [README.md](#) (File gives a summary of the project, its files, and how to run it)
- [design.pdf](#) (File gives a more in depth description of the project, its files, and how to run it)

2.2 Program Options

Note: By default, both decode and encode programs take input from stdin and output to stdout.

The encode program supports the following options:

- -i <input file> : specify input to compress.
- -o <output file>: specify output file to decompress.
- -v: display compression stats.

The decode program supports the following options:

- -i <input file> : specify input to decompress.
- -o <output file>: specify output file of decompressed input.
- -v: display decompression stats.

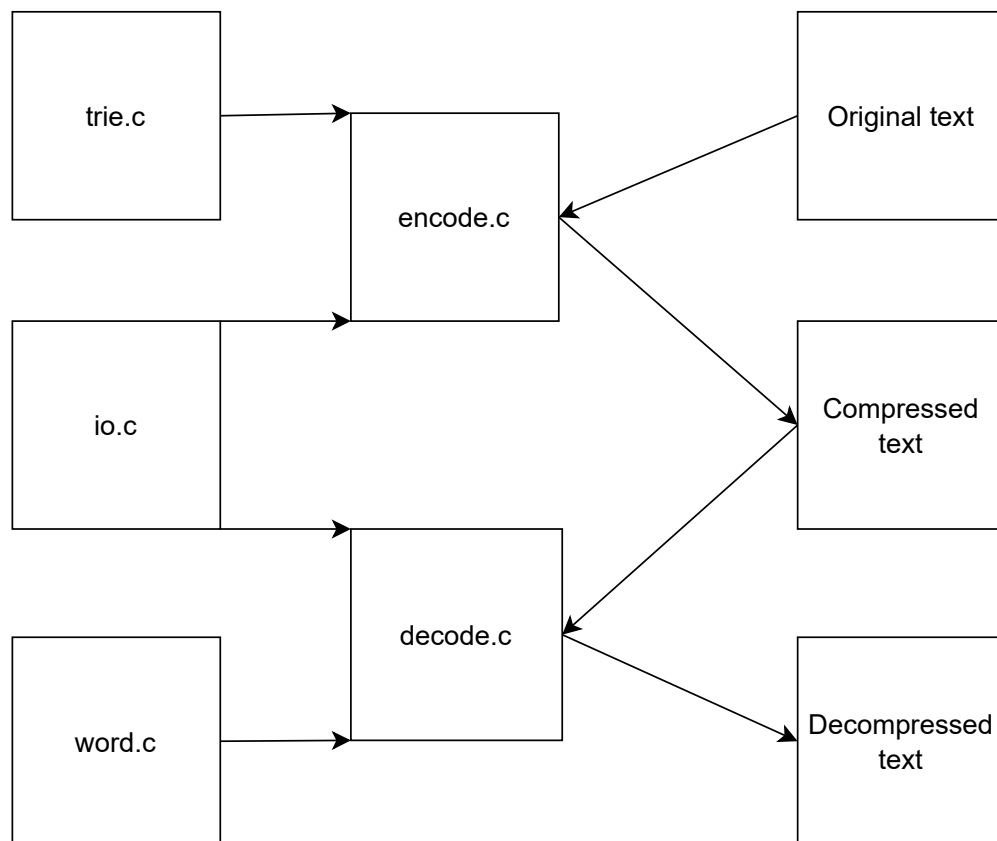


Figure 1: The block diagram above depicts the flow of how the program files interact with one another and with the file to compress. As shown, `encode.c` uses `trie.c` and `io.c` to compress the Original file and outputs the Compressed file. Similarly, `decode.c` uses `word.c` and `io.c` to decompress the Compressed file and outputs the Decompressed file. Made with draw.io [?].

3 Running The Program

To run the program, the encode and decode executables must first be created. To make the executables type "make" from the same directory as the Makefile. A data can then be compressed using the command "./encode". The -i |infile| and the -o |outfile| can be used to specify a file to compress and where to output the compressed data respectively. The -v option can be used to print out stats pertaining the encoding. Decoding a compressed file can be done using the command "./decode". The -i |infile| and the -o |outfile| can be used to specify the compressed file and the location to output the decompressed output respectively. Similarly to encode, the -v option can be used to output decompression stats. Encode and decode binary executables can then be removed by typing "make clean"

Summarized Instructions (options in "[]" are optional):

1. From the same directory as the Makefile type "make"
2. Encode: "./encode [-i |infile|] [-o |outfile|] [-v]"
3. Decode: "./decode [-i |infile|] [-o |outfile|] [-v]"
4. Clean: "make clean"

4 Compression/Decompression

4.1 TRIE - trie.c

The trie.c file contains functions used to create and manipulate a trie abstract data structure. A trie is a prefix tree, where each node contains n child nodes, 256 for this project. The trie (and its related functions) are used for the compression of files, since it is an optimal way of check for existing prefixes, or words.

4.1.1 trie.c Pseudocode

```
1 TrieNode data structure {
2     Pointer to an array of 256 TrieNodes (TrieNode *children[256])
3     code
4 }
```

```
1 trie_node_create(code){
2     Allocate space the size of a trienode to a new trie node.
3     for the children of the trie node (256 indexes) {
4         trie->children[index] = NULL
5     }
6     trie_node->code = code
7     return trie_node
8 }
```

```

1 trie_node_delete(trie_node){
2     free memory allocated to trie_node
3     return
4 }

```

```

1 trie_create(){
2     Create a new trie node using trie_node_create with EMPTY_CODE as the input argument
3     return a pointer to the created trie node
4 }

```

```

1 trie_reset(root){
2     for the children of the root (256 indexes) {
3         if the child node is not NULL{
4             trie_reset(root->child_node[index]) // Recursive call using child node as input
5             Free the node using trie_node_delete(root->child_node[index])
6             Set the child node to NULL
7         }
8     }
9     return
10 }

```

```

1 trie_delete(root) {
2     for the children of the root (256 indexes) {
3         if the child is not NULL{
4             trie_delete(root->child_node[index]) // Recursive call using child node as input
5         }
6     }
7     Free the node here to ensure all nodes get freed including the root
8     return
9 }

```

```

1 trie_step(trie_node, sym){
2     if the trie_node is not NULL{
3         return trie_node's child at index sym (trie_node->child[sym])
4     }
5     Otherwise return NULL
6 }

```

4.2 WORD - word.c

The word.c file contains functions used to create and manipulate Word and WordTable data structures. A Word contains both a symbol string (words/prefixes) and a code corresponding to the string. A WordTable is an array of Words that starts with a Word containing an empty symbol string and START_CODE; it ultimately acts as a table of words and codes. Although a trie could have been used, word.c is used for decompression (decode.c) as it makes looking up code to word translations quicker, thus significantly reducing the time it takes to decompress a file.

4.2.1 word.c Pseudocode

```

1 Word data structure {
2     character pointer syms
3     length of syms
4 }

```

```

1 word_create(syms, len){
2     Allocate space the size of a Word to a new_Word.
3     if len is not 0{
4         allocate len amount of uint8_t's to the new_Words syms attribute (new_Word->syms)
5     }
6     if len is 0 then allocate 1 uint8_t to the new_Words syms attribute
7     copy the input syms to the new_Word syms attribute (copy syms to new_Word->syms)
8     return new_Word
9 }

```

```

1 word_append_sym(word, sym){
2     Create a copy_word of the input word using word_create and passing in words attributes as
3     arguments (word_create(word->syms, word->len))
4
5     if the syms attribute of word is not empty{
6         reallocate memory to the syms attribute of copy_word by 1 uint8_t to make room for
7         the sym to appended.
8     }
9     // No need to reallocate memory if syms attribute of word is empty since it was allocated
10    // 1 uint8_t when it was created.
11    increase the len attribute of copy_word by 1
12    add sym to the last index of copy_word's attribute syms.
13    return copy_word
14 }

```

```

1 word_delete(word){
2     if words attribute syms is not pointing to NULL free syms
3     free word
4     return
5 }

```

```

1 wt_create(){
2     Allocate memory for a word_table of (UINT16_MAX - 1) Words
3     Create a word with empty sym ("") and len 0 and place it in word_table at index
4     EMPTY_CODE
5
6     return word_table
7 }

```

```

1 wt_reset(word_table){
2     for all the words in the word table (UINT16_MAX - 1){
3         if the Word is not NULL (if it exists) {
4             Delete the word (word_delete())
5             Set the word to NULL
6         }
7     }
8     return
9 }

```

```

1 wt_delete(word_table){
2     for all the words in the word table (UINT16_MAX - 1){
3         if the word is not null delete the word (word_delete())
4     }
5     free memory allocated to the word table
6     return
7 }

```

4.3 IO - io.c

The io.c file contains functions that are used for both compressing and decompressing files (encode.c and decode.c) to read and write from and to files. Reading and writing is done in blocks of 4KB to make the I/O more efficient. When reading, a new block of data is read in only when the previously read in data has been fully processed. Similarly, when writing, the buffered data is only written when the buffer is full.

4.3.1 io.c Pseudocode

```

1 struct FileHeader {
2     uint32_t magic
3     uint16_t protection
4 }

```

```

1 read_header(infile, header){
2     read a block the size of a FileHeader into the input FileHeader "header"
3     if the order of the bytes on the system is big endian{
4         swap the protection attribute of the input header using (swap16(header->protection))
5         swap the magic attribute of the input header using (swap32(header->magic))
6     }
7     return;
8 }

```

```

1 write_header(outfile, header){
2     if the order of the bytes on the system is big endian{
3         swap the protection attribute of the input header using (swap16(header->protection))
4         swap the magic attribute of the input header using (swap32(header->magic))
5     }
6     write header to the outfile
7     return;
8 }

```

```

1 read_sym(infile, sym){
2     read a block of data from the input file to a buffer if the end of the buffer was reached
3     or if we have never read into the buffer. Reset the index of the buffer.
4
5     copy the symbol (byte) at the current index of the buffer into input sym
6     increase the index of the buffer
7
8     if sym is 0 then return false
9     otherwise return true
10 }

```

```

1 buffer_pair(outfile, code, sym, bitlen){
2     for (int i = 0; i < bitlen; i++ (length of code)){
3         if the bit_index/8 is at the buffer limit{
4             write the buffer out to the output file
5             reset the buffer
6             reset the bit_index
7         }
8         check if the bit of the input code at index i is set.
9         if the bit is set{
10             set the bit of the buffer at index (buffer[bit_index / 8] since bit_index is
11             kept track of in bits not bytes)
12         }
13         increase the bit_index value by 1
14     }
15
16     for (int j = 0; j < CHAR_BIT; j++){
17         if the bit_index/8 is at the buffer limit{
18             write the buffer out to the output file
19             reset the buffer
20             reset the bit_index
21         }
22         check if the bit of the input sym at index j is set.
23         if the bit is set{
24             set the bit of the buffer at index (buffer[bit_index / 8] since bit_index is
25             kept track of in bits not bytes)
26         }
27         increase the bit_index value by 1
28     }
29     return
30 }

```

```

1 flush_pairs(outfile){
2     write the pair buffer to the outfile
3     reset the buffer
4     return
5 }

```

```

1 read_pair(infile, code, sym, bitlen){
2     for (i = 0; i < bitlen; i++) {
3         if the bit_index/8 is at the buffer limit or no blocks have been read{
4             reset the buffer
5             read in a new block from infile into the buffer
6             reset the bit_index
7         }
8         check if the bit of the buffer[bit_index/8] at index i is set.
9         if the bit is set{
10            set the bit of input code at index i
11        }
12        increase the bit_index by 1
13    }
14
15    for (j = 0; j < CHAR_BITS; j++) {
16        if the bit_index/8 is at the buffer limit or no blocks have been read{
17            reset the buffer
18            read in a new block from infile into the buffer
19            reset the bit_index
20        }
21        check if the bit of the buffer[bit_index/8] at index j is set.
22        if the bit is set{
23            set the bit of input sym at index j
24        }
25        increase the bit_index by 1
26    }
27
28    if code is the STOP_CODE return false
29    otherwise return true
30 }

```

```

1 buffer_word(outfile, word){
2     for all the character of the word attribute sym {
3         for (j = 0; j < CHAR_BITS; j++ (all bits in character)) {
4             check if the bit of the character is set
5             if it is set{
6                 set the bit of buffer[bit_index/8] at index j
7             }
8         }
9         increase the bit_index by one
10        if the bit_index/8 reaches the buffer limit{
11            write the buffer to the outfile
12            reset the buffer
13            reset the bit_index
14        }
15    }
16    return
17 }

```



```

1 flush_word(outfile) {
2     write the word buffer to the outfile
3     reset the word buffer
4     return
5 }

```

```

1 read_buffer(buffer, byte_num, infile){
2     reset the buffer
3     read byte_num bytes into the buffer from infile
4     return
5 }

```

```

1 write_buffer(buffer, byte_num, outfile){
2     write byte_num bytes of buffer into the outfile
3     reset the buffer
4     return
5 }

```

4.4 Encode - encode.c

"encode.c" is the primary file responsible for compressing files. It contains the core compression algorithm, leveraging "io.c" and "trie.c". Additionally, "encode.c" handles program options as well as file management.

4.4.1 encode.c Pseudocode

```

1 bitlen(x){
2     if x is not 0{
3         while x is 1{
4             add one to count
5             shift x to the left by 1
6         }
7         return count
8     }
9     else{
10        return 1
11    }
12 }

```

```

1 encode_header(outfile, infile){
2     Allocate memory of one FileHeader to a new FileHeader.
3     Add the magic attribute to the new FileHeader (new_FH->MAGIC)
4     Add the protection header to the new FileHeader (permissions of the infile)
5     write the header to the outfile
6     free the memory allocated to the new header
7     return
8 }

```

```

1 compress_file(root, outfile, infile){
2     set curr_node to the root
3     set curr_sym to 0
4     while there are symbols in the infile{
5         set next node to the output of trie_step(curr_node, curr_sym)
6         if next node is NULL{
7             set prev_node to curr_node and curr_node to next_node
8         }
9         else{
10            write the pair (code and sym) to the outfile using buffer_pair()
11            Add a new child for curr_node at index curr_sym using next_code
12            make the curr_node the root
13            increase next_code by 1
14        }
15        if next code has reached its limit (MAX_CODE){
16            reset the trie (trie_rest)
17            set curr_node to the root
18            set the next_code to the START_CODE
19        }
20        set prev_sym to curr_sym
21    }
22
23    if curr_node is not the root{
24        output the last pair in the file
25        increase next_code by 1
26    }
27    Output the STOP_CODE to the outfile
28    flush the pair buffer
29 }

```

4.5 Decode - decode.c

"decode.c" is the primary file responsible for decompressing files. It contains the core decompression algorithm, leveraging "io.c" and "word.c". Additionally, "decode.c" handles program options as well as file management.

4.5.1 decode.c Pseudocode

```

1 bitlen(x){
2     if x is not 0{
3         while x is 1{
4             add one to count
5             shift x to the left by 1
6         }
7         return count
8     }
9     else{
10        return 1
11    }
12 }

```

```

1 decode_header(infile){
2     Allocate memory of one FileHeader to a new FileHeader.
3     read the header from the infile to the new FileHeader.
4
5     if the header magic does not match the expected MAGIC{
6         free the new header
7         print an error message and exit the program
8     }
9     free the header
10    return
11 }

```

```

1 decompress_file(table, outfile, infile){
2     set next_code to START_CODE
3     while there are pairs in the infile{
4         Add the appended symbol/word to the table at index next_code
5         increase next_code by 1
6         if next code has reached the code limit (MAX_CODE){
7             reset the table
8             set the next_code to the START_CODE
9         }
10    }
11    flush the word buffer to the outfile
12    return
13 }

```

4.6 Endian - endian.h

Code provided by Professor Darrell Long. This file contains functions to check if the order of bytes of a system is big or little endian, as well as functions for swaping the endianness of uint16_t, uint32_t, and uint64_t.

4.6.1 endian.h Code

```

1 is_big(){
2     union{
3         uint8_t bytes[2];
4         uint16_t word;
5     }test;
6     test.word = 0xFF00;
7     return test.bytes[0];
8 }

```

```

1 is_little(){
2     return ! is_big();
3 }

```

```

1 static inline uint16_t swap16(uint16_t x){
2     uint16_t result = 0;
3     result |= (x & 0x000000FF) << 24;
4     result |= (x & 0x0000FF00) << 8;
5     result |= (x & 0x00FF0000) >> 8;
6     result |= (x & 0xFF000000) >> 24;
7     return result;
8 }

```

```

1 static inline uint32_t swap32 (uint32_t x){
2     uint32_t result = 0;
3     result |= (x & 0x000000FF) << 24;
4     result |= (x & 0x0000FF00) << 8;
5     result |= (x & 0x00FF0000) >> 8;
6     result |= (x & 0xFF000000) >> 24;
7     return result;
8 }

```

```

1 static inline uint64_t swap64(uint64_t x){
2     uint64_t result = 0;
3     result |= (x & 0x00000000000000FF) << 56;
4     result |= (x & 0x000000000000FF00) << 40;
5     result |= (x & 0x0000000000FF0000) << 24;
6     result |= (x & 0x00000000FF000000) << 8;
7     result |= (x & 0x000000FF00000000) >> 8;
8     result |= (x & 0x0000FF0000000000) >> 24;
9     result |= (x & 0x00FF000000000000) >> 40;
10    result |= (x & 0xFF00000000000000) >> 56;
11    return result;
12 }

```

4.7 Code - code.h

This file contains important codes that are used for the compression and decompression of files (word.c and trie.c).

Codes defined in code.h:

- STOP_CODE 0
- EMPTY_CODE 1
- START_CODE 2
- MAX_CODE UINT16_MAX

5 Testing

Note: The tests can be ran all at once using the test.sh script using the command **”./test.sh”** from the same directory as the test.sh file.

- test.sh (Main test file that can be used to run all other tests)

- `small_text.sh` (Test verifying the compression and decompression of a small text file)
- `large_text.sh` (Test verifying the compression and decompression of a the bible as a text file)
- `small_binary.sh` (Test verifying the compression and decompression of a small png file)
- `large_binary.sh` (Test verifying the compression and decompression of a large png file)
- `valgrind_text.sh` (Test checking the program for memory leaks while compressing and decompressing the bible as a text file)
- `utils.sh` (File that contains functions used in `test.sh` (to clean temporary files and check if files exist))

6 Other Files

6.1 Makefile

The Makefile is a file written in make file language and is used to create an encode and decode binary executable files. The executables are created using the command 'make' and removed using 'make clean', both must be run in the same directory as the Makefile.

6.2 README.md

The README.md file gives a brief overview of the program, its files, and how to run it.

6.3 design.pdf

The design.pdf is this file. It gives a more in depth explanation about the program, its files, and how to run it.