

Project 5: Parallel Zip

Daniil Komov, Igor Zimarev

Link to source code Github repository: https://github.com/lkomovONE/Parallel_Zip

For this mini-project, we made two C programs that compress/decompress files in RLE encoding using parallel multi-threading.

Program pzip.c takes in one or multiple files (works best with .txt files), compresses the file/files using parallel threads and outputs the compressed file.

Program punzip.c takes in one or multiple compressed files (works best with compressed .txt files), decompresses the file/files using parallel threads and outputs the decompressed file.

Files pzip.c and punzip.c are the files containing the code for both programs. Both the programs can be compiled with gcc using this command: “gcc -o pzip pzip.c” or “gcc -o punzip punzip.c”

The program has been designed to work primarily with .txt files, and is limited in a way, but it makes sense to test it with other different file extensions (for example to compress .c file)

The repository already includes both of the compiled programs. There are several ways to run the programs:

-1 argument:

```
prompt> ./pzip inputFile > outputFile.z
```

This prompt will make program use 1 input file to compress it into “outputFile.z”. It’s important to note that compressed file has to be specified, and it has to be with .z extension.

-Many arguments:

```
prompt> ./pzip inputFile1 inputFile2 ... > outputFile.z
```

This prompt will make program take few input files, and compress them into one output file. It’s important to note that while there may be several input files, they are all being written into 1 compressed output file with .z extension.

-punzip program:

```
prompt> ./punzip compressed_inputFile > decompressed_outputFile.txt
```

or

```
prompt> ./punzip compressed_inputFile compressed_inputFile2 >  
decompressed_outputFile.txt
```

Similar execution is for punzip program: there may be several files to decompress, but they are all going into one unified decompressed file. The input files have to be accessible by the program.

Error handling

Before using the file to compress/decompress, make sure you have file permissions set-up as Read & Write for the system.

- If the file is corrupt or unavailable, error will pop up “error: cannot open file ‘<file>”.
- In case of internal memory allocation issue, an error message will appear “malloc failed”.

The program has no file size limitations, however, very large files may affect system performance or compression may fail due to system memory constraints.

Considerations

The project codebase has been developed based on different considerations:

Parallelizing the compression.

The program simply divides the files into several (as many as the number of threads) segments, and each segment is being treated by assigned to it thread. Each segment gets only 1 thread, and all segments are made equal since the segment size is determined as file size if divided by the number of threads. This way system handles the file using several threads at once.

Determining how many threads to create

The program reads the number of system processors available, then accordingly creates as many threads as the number of these processors.

Efficiently performing each piece of work.

All segments are equally divided, and simple command in a loop is used to encode the contents of each segment. Memory mapping is used to access the files and parallelization is included. Thus, program can be treated as efficiency-oriented.

Accessing the input file efficiently

To access input files the program simply opens the file, then uses `mmap()` command to map the file into the memory address space.

Use cases examples

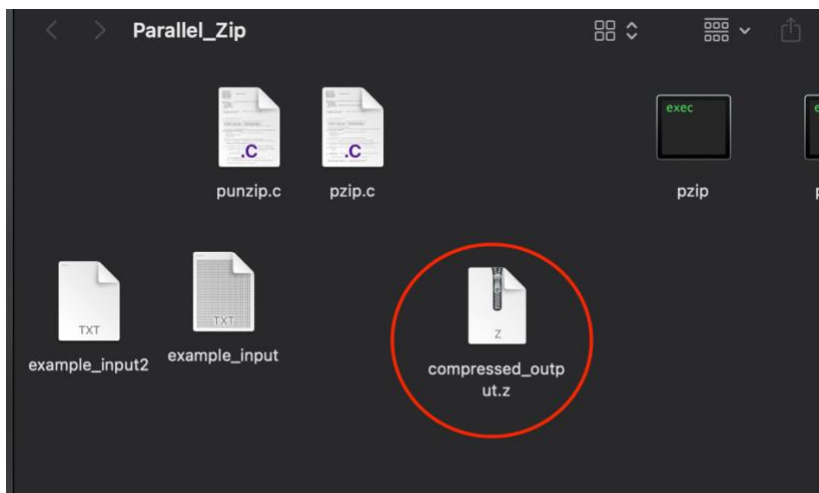
Repository contains example input files that are used here.

-1 argument:

```
prompt> ./pzip example_input.txt > compressed_output.z
```

```
@Daniils-MacBook-Air Parallel_Zip % ./pzip example_input.txt > compressed_output.z
@Daniils-MacBook-Air Parallel_Zip %
```

The program has no output in terminal, however, it creates the file inside the directory:



-Many arguments:

```
prompt> ./pzip example_input.txt example_input2.txt> compressed_output.z
```

The program never prints anything in terminal if the execution is successful. In this scenario the output of the program will be exactly same as in previous use case, except that the compressed_output.z file will have different contents.

-punzip:

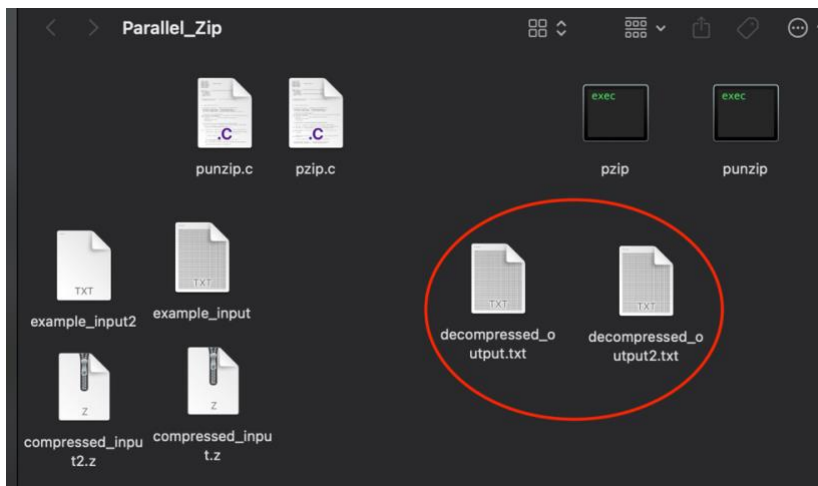
```
prompt> ./punzip compressed_input.z > decompressed_output.txt
```

or

```
prompt> ./punzip compressed_input.z compressed_input2.z > decompressed_output.txt
```

```
ikomovone@Daniils-MacBook-Air Parallel_Zip % ./punzip compressed_input.z > decompressed_output.txt
ikomovone@Daniils-MacBook-Air Parallel_Zip % ./punzip compressed_input.z compressed_input2.z > decompressed_output2.txt
ikomovone@Daniils-MacBook-Air Parallel_Zip %
```

The program executes the same way as pzip, except it will add decompressed file/files in the directory that would contain the original text:



Screenshots of the punzip.c code:

```
C punzip.c M X
C punzip.c > main(int, char * [])
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <fcntl.h>
4  #include <string.h>          //importing libraries
5  #include <sys/stat.h>
6  #include <sys/mman.h>
7  #include <unistd.h>
8  #include <pthread.h>
9  #include <ctype.h>
10
11
12
13  typedef struct {              //initializing certain data structure for each thread, with start of data segment, its length, and thread index
14      char *seg_start;
15      size_t seg_length;
16      size_t th_index;
17      char *output_buffer;
18      size_t output_length;
19  } Th_data;                    //Reference: https://www.educative.io/answers/how-to-use-the-typedef-struct-in-c?t
20
21
22
23  void message_printer(const char *msg) {    //Function for printing errors and messages
24      fprintf(stderr, "%s\n", msg);
25  }
26
27
28
29
30
31  //-----Segment decompression function-----
32
33
34  void *seg_decompression(void *th_arg) {    //Establishing function for segment decompression
35
```

```
C punzip.c M X
C punzip.c > main(int, char * [])
31  //-----Segment decompression function-----
32
33
34  void *seg_decompression(void *th_arg) {    //Establishing function for segment decompression
35
36      Th_data *seg_data = (Th_data *)th_arg; //establishing thread data instance
37      char *seg_start = seg_data->seg_start; //initializing segment starting point
38      size_t seg_length = seg_data->seg_length; //initializing length of the segment
39      size_t th_index = seg_data->th_index; //Initializing thread index value for possible troubleshooting
40
41
42      //creating allocation for decompressed segment data
43
44      char *decomp_output = malloc(seg_length * 10); // Allocating memory for the output, multiplying the length of the segment so that there is always
45
46      if (!decomp_output) {                    //Adding error handling in case allocation fails
47          message_printer("malloc failed");
48
49          pthread_exit(NULL);
50      }
51
52      size_t decomp_buff_position = 0; //initializing variable for positioning in the buffer
53
54      for (size_t i = 0; i < seg_length; i) { //Using for-loop to go over each encoded character inside the segment
55
56          if (i + sizeof(int) + 1 > seg_length) { //Print error message if there's space for the 4-byte count reading
57
58              message_printer("Error, the file seems to be corrupted. The size for the count encoding doesn't match!");
59
60              free(decomp_output);
61
62              pthread_exit(NULL);
63          }
64      }
65
```

```

34 void *seg_decompression(void *th_arg) { //Establishing function for segment decompression
55     for (size_t i = 0; i < seg_length; i) { //Using for-loop to go over each encoded character inside the segment
66
67         int c_count; //Initializing count value for RLE encoding.
68
69         //Reading the count from the pair
70
71         memcpy(&c_count, seg_start + i, sizeof(int)); //Reference: https://www.javatpoint.com/memcpy-in-c?t
72
73         i += sizeof(int); // Moving the index to the character
74
75         // Reading the character value by initializing the char variable
76
77         char enc_character = *(seg_start + i);
78
79         i++; // Moving the index to the next pair
80
81         // Checking if the count value is correct
82         if (c_count < 0 || decomp_buff_position + c_count > seg_length * 10) {
83
84             message_printer("Error, count value is incorrect. Corrupted File!");
85
86             free(decomp_output);
87
88             pthread_exit(NULL);
89         }
90
91         // Using for-loop to write the character into the buffer x number of times
92         for (int j = 0; j < c_count; j++) {
93
94             decomp_output[decomp_buff_position++] = enc_character;
95
96         }
97     }

```

```

96     }
97 }
98
99
100 seg_data->output_buffer = decomp_output; //Writing the result into the buffer
101 seg_data->output_length = decomp_buff_position; //Updating the buffer position
102
103
104
105 pthread_exit(NULL); //exiting thread
106 }
107
108
109
110
111
112
113
114
115
116
117 //-----Main function-----
118
119
120
121
122
123
124 int main(int arg_counter, char *arg_select[]) { //establishing the main function
125
126

```

```

22
23
24 int main(int arg_counter, char *arg_select[]) { //establishing the main function
25
26
27
28     if (arg_counter < 2) { //Handling incorrect usage cases by throwing usage suggestion
29
30         message_printer("usage: punzip <file1> <file2> ... > <outputFile1> <outputFile2>"); //calling message printer to print correct usage
31
32
33         exit(1); //exiting the program
34     }
35
36
37
38
39     //Next step is to establish the threads, doing it based on the number of processors system has
40
41     int th_number = sysconf(_SC_NPROCESSORS_ONLN); //getting number of threads by getting the number of processors of the system. Using sysconf
42
43     //Reference for the sysconf(): https://man7.org/linux/man-pages/man3/sysconf.3.html
44
45
46
47     if (th_number < 1) { //in case sysconf fails, we set the number of threads to 1 as default number (1 thread)
48         th_number = 1;
49     }
50
51
52     pthread_t threads[th_number]; //Initializing the threads based on the established amount
53
54     Th_data th_data[th_number]; //Establishing Th_data instances (amount of instances based on thread number). Th_data data structure has been s
55

```

```

8
9
10 for (int i = 1; i < arg_counter; i++) { //establishing a loop for all the arguments (input files)
11
12     int input_file = open(arg_select[i], O_RDONLY); //opening the file (one by one as it's a loop)
13
14
15     if (input_file < 0) { //Error handling for file opening
16
17         message_printer("Could not open the file");
18
19         exit(1);
20     }
21
22
23     struct stat file_stats; //initializing instance for collecting file statistics (for example size)
24
25     //Reference for the stat: http://codewiki.wikidot.com/c:system-calls:fstat?t
26
27
28
29     if (fstat(input_file, &file_stats) != 0) { //Error handling for file statistics
30
31         message_printer("Couldn't get the size of the file");
32
33         close(input_file);
34
35         exit(1);
36     }
37
38
39
40     char *data = mmap(NULL, file_stats.st_size, PROT_READ, MAP_PRIVATE, input_file, 0); //mapping file data into memory for easy access
41

```

```

if (data == MAP_FAILED) {    //Error handling for file data mapping

    message_printer("Couldn't map the file data to memory");

    close(input_file);

    exit(1);
}

size_t seg_size = file_stats.st_size / th_number; //establishing segments' sizes based on file size divided by number of threads, so that

for (int k = 0; k < th_number; k++) {

    size_t seg_start_offset = k * seg_size; //Getting start offset for the thread to know where to start
    size_t seg_end = (k == th_number - 1) ? file_stats.st_size : seg_start_offset + seg_size; //Determining the end of the segment

    if (k > 0) {            //Using if-conditions and a loop to adjust start offset for the segments that are following after the 1st one
        while (seg_start_offset < seg_end && (seg_start_offset % (sizeof(int) + 1) != 0)) {
            seg_start_offset++;
        }
    }

    while (seg_end < file_stats.st_size && (seg_end % (sizeof(int) + 1) != 0)) { //Using a loop and if condition to adjust the end
        seg_end++;
    }

    if (seg_start_offset >= file_stats.st_size || seg_end > file_stats.st_size) { //Adding error handling for segment boundaries
        message_printer("Out of bounds error in segment. The compressed file may be corrupt");
        exit(1);
    }

    //Establishing thread data for creating the thread
    th_data[k].seg_start = data + seg_start_offset;
    th_data[k].seg_length = seg_end - seg_start_offset;
    th_data[k].th_index = k;

    //printf("Thread %d: seg_start_offset = %zu, seg_end = %zu, seg_length = %zu\n", k, seg_start_offset, seg_end, th_data[k].seg_length);

    pthread_create(&threads[k], NULL, seg_decompression, &th_data[k]); //Creating thread
}

for (int j = 0; j < th_number; j++) { //Using for-loop for joining all the threads

    pthread_join(threads[j], NULL);

    fwrite(th_data[j].output_buffer, sizeof(char), th_data[j].output_length, stdout); //Writing to a file
}

```



```

250     for (int j = 0; j < th_number; j++) { //Using for-loop for joining all the threads
251
252
253         pthread_join(threads[j], NULL);
254
255         fwrite(th_data[j].output_buffer, sizeof(char), th_data[j].output_length, stdout); //Writing to a file
256
257         free(th_data[j].output_buffer); //Freeing the output buffer
258     }
259
260
261
262
263     munmap(data, file_stats.st_size); //Freeing mapped memory
264
265     close(input_file); //closing input file
266 }
267
268
269
270
271
272     return 0;
273
274 }
275
276

```

Screenshots of the pzip.c code:

```
C pzip.c M X
C pzip.c > main(int, char* [])
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <fcntl.h>
4  #include <string.h>           //importing libraries
5  #include <sys/stat.h>
6  #include <sys/mman.h>
7  #include <unistd.h>
8  #include <pthread.h>
9
10
11
12  typedef struct {              //initializing certain data structure for each thread, with start of data segment, its length, and thread index
13      char *seg_start;
14      size_t seg_length;
15      size_t th_index;
16      char *output_buffer;     //Establishing output buffer and length for smooth result printing
17      size_t output_length;
18  } Th_data;                    //Reference: https://www.educative.io/answers/how-to-use-the-typedef-struct-in-c?t
19
20
21
22
23  void message_printer(const char *msg) {    //Function for printing errors and messages
24      fprintf(stderr, "%s\n", msg);
25  }
26
27
28
29
30
31  //-----Segment compression function-----
32
33
34  void *seg_compression(void *th_arg) {      //Establishing function for segment compression
35
36  C pzip.c > main(int, char* [])
37      void *seg_compression(void *th_arg) {    //Establishing function for segment compression
38
39          Th_data *seg_data = (Th_data *)th_arg; //establishing thread data instance
40          char *seg_start = seg_data->seg_start; //initializing segment starting point
41          size_t seg_length = seg_data->seg_length; //initializing length of the segment
42
43          //creating allocation for compressed segment data
44
45          char *comp_output = malloc(seg_length * 5); // Allocating memory for the output, multiplying the length of the segment so that there is always
46
47          if (!comp_output) {                    //Adding error handling in case allocation fails
48              message_printer("malloc failed");
49              pthread_exit(NULL);
50          }
51
52          size_t comp_buff_position = 0; //initializing variable for positioning in the buffer
53
54          for (size_t i = 0; i < seg_length; i) { //Using for-loop to go over each character inside the segment
55
56              char character = seg_start[i]; //initializing the current character, setting it to the starting point in the segment
57              size_t counter = 1; //Initializing counter for RLE encoding. Counts the same character trend.
58
59              while (i + counter < seg_length && seg_start[i + counter] == character) { //Using while-loop to count characters of same type
60
61                  counter++;
62              }
63
64              //writing the data to the memory buffer
65
66
67
68
```

```

55     for (size_t i = 0; i < seg_length; i++) { //Using for-loop to go over each character inside the segment
56
57         //First writing the binary count
58         memcpy(comp_output + comp_buff_position, &counter, sizeof(int)); //Reference: https://www.javatpoint.com/memcpy-in-c?t
59
60         comp_buff_position += sizeof(int);
61
62         memcpy(comp_output + comp_buff_position, &character, sizeof(char)); //Then writing the character itself
63
64         comp_buff_position += sizeof(char);
65
66         i += counter; // Adjusting i variable for the loop to progress through the characters
67     }
68
69     seg_data->output_buffer = comp_output; //Writing the output data to output buffer
70     seg_data->output_length = comp_buff_position; //Adjusting the position of buffer
71
72     pthread_exit(NULL); //exiting thread
73 }
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97

```

```

98
99
100
101
102
103
104 //-----Main function-----
105
106
107
108
109 int main(int arg_counter, char *arg_select[]) { //establishing the main function
110
111     if (arg_counter < 2) { //Handling incorrect usage cases by throwing usage suggestion
112         message_printer("usage: pzip <file1> <file2> ... > <outputFile>"); //calling message printer to print correct usage
113
114         exit(1); //exiting the program
115     }
116
117
118
119
120
121
122
123
124 //Next step is to establish the threads, doing it based on the number of processors system has
125
126 int th_number = sysconf(_SC_NPROCESSORS_ONLN); //getting number of threads by getting the number of processors of the system. Using sysconf f
127
128 //Reference for the sysconf(): https://man7.org/linux/man-pages/man3/sysconf.3.html
129
130
131
132 if (th_number < 1) { //in case sysconf fails, we set the number of threads to 1 as default number (1 thread)
133     th_number = 1;
134 }
135

```

```

135
136
137 pthread_t threads[th_number]; //Initializing the threads based on the established amount
138
139 Th_data th_data[th_number]; //Establishing Th_data instances (amount of instances based on thread number). Th_data data structure has been i
140
141
142
143
144 for (int i = 1; i < arg_counter; i++) { //establishing a loop for all the arguments (input files)
145
146     int input_file = open(arg_select[i], O_RDONLY); //opening the file (one by one as it's a loop)
147
148
149     if (input_file < 0) { //Error handling for file opening
150
151         message_printer("Could not open the file");
152
153         exit(1);
154     }
155
156
157     struct stat file_stats; //initializing instance for collecting file statistics (for example size)
158
159     //Reference for the stat: http://codewiki.wikidot.com/c:system-calls:fstat?t
160
161
162
163     if (fstat(input_file, &file_stats) != 0) { //Error handling for file statistics
164
165         message_printer("Couldn't get the size of the file");
166
167         close(input_file);
168
169         exit(1);
170     }
171
172
173
174     char *data = mmap(NULL, file_stats.st_size, PROT_READ, MAP_PRIVATE, input_file, 0); //mapping file data into memory for easy access
175
176     if (data == MAP_FAILED) { //Error handling for file data mapping
177
178         message_printer("Couldn't map the file data to memory");
179
180         close(input_file);
181
182         exit(1);
183     }
184
185
186
187     size_t seg_size = file_stats.st_size / th_number; //establishing segments' sizes based on file size divided by number of threads, so that
188
189
190
191     for (int k = 0; k < th_number; k++) { //Using for-loop to create each thread 1 by 1, for file compression
192
193         //Difining each element of the thread's data structure, such as start of segment, segment length and thread index
194
195
196         th_data[k].seg_start = data + k * seg_size; //Difining start of segment. The calculation is an offset that the thread has to make to
197
198         th_data[k].seg_length = (k == th_number - 1) ? (file_stats.st_size - k * seg_size) : seg_size; //Difining length of segment, includes
199

```

```

197 //for (int k = 0; k < th_number; k++) { //Using for-loop to create each thread & seg for file compression
198     th_data[k].seg_length = (k == th_number - 1) ? (file_stats.st_size - k * seg_size) : seg_size; //Defining length of segment, includes \
199     th_data[k].th_index = k; //Setting index of thread as k number of the for-loop (simple counting)
200
201
202
203     pthread_create(&threads[k], NULL, seg_compression, &th_data[k]); //creating thread, function for compressing the segment specified
204
205
206 }
207
208
209
210
211 for (int j = 0; j < th_number; j++) { //Using for-loop for joining all the threads
212
213     pthread_join(threads[j], NULL);
214
215     fwrite(th_data[j].output_buffer, sizeof(char), th_data[j].output_length, stdout); //Writing to the file
216
217
218     free(th_data[j].output_buffer); //Freeing the output buffer
219
220 }
221
222
223
224
225 munmap(data, file_stats.st_size); //Freeing mapped memory
226
227 close(input_file); //closing input file
228 }

```

```

207
208
209
210
211 for (int j = 0; j < th_number; j++) { //Using for-loop for joining all the threads
212
213     pthread_join(threads[j], NULL);
214
215     fwrite(th_data[j].output_buffer, sizeof(char), th_data[j].output_length, stdout); //Writing to the file
216
217
218     free(th_data[j].output_buffer); //Freeing the output buffer
219
220 }
221
222
223
224
225 munmap(data, file_stats.st_size); //Freeing mapped memory
226
227 close(input_file); //closing input file
228 }
229
230
231
232
233 return 0;
234
235 }
236
237

```