

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

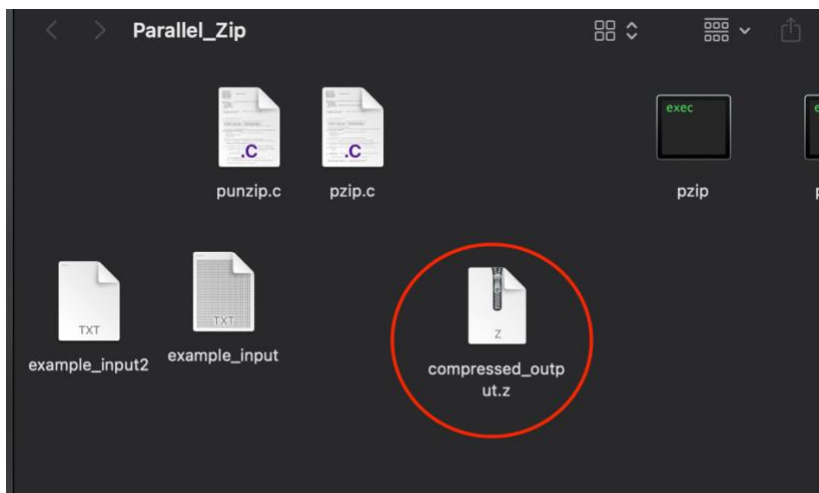
Repository contains example input files that are used here.

-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:

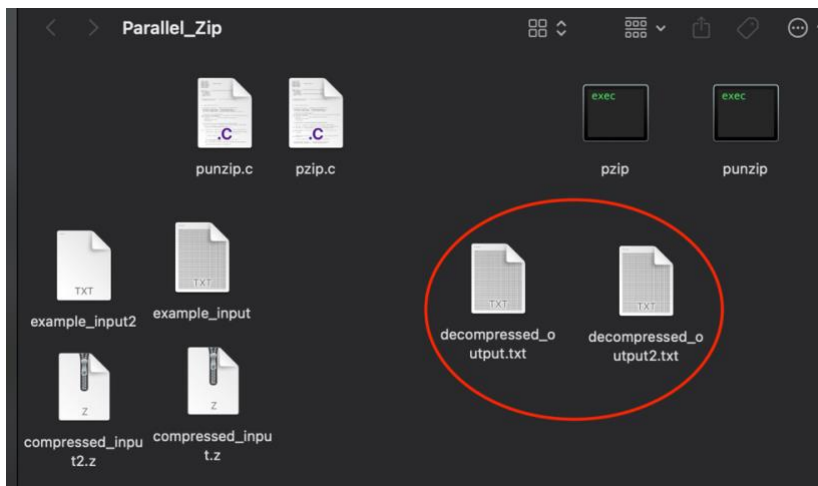
```
prompt> ./punzip compressed_input.z > decompressed_output.txt
```

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 pzip.c code:

```

C pzip.c > ...
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;
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

```

```

33
34  void *seg_compression(void *th_arg) {     //Establishing function for segment compression
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
40
41      //creating allocation for compressed segment data
42
43      char *comp_output = malloc(seg_length * 2); // Allocating memory for the output, doubling the length of the segment so that there is always enough
44
45      if (!comp_output) {                //Adding error handling in case allocation fails
46          message_printer("malloc failed");
47          pthread_exit(NULL);
48      }
49
50
51      size_t comp_buff_position = 0; //initializing variable for positioning in the buffer
52
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
60          while (i + counter < seg_length && seg_start[i + counter] == character) { //Using while-loop to count characters of same type
61              counter++;
62          }
63
64          comp_buff_position += sprintf(comp_output + comp_buff_position, "%zu%c", counter, character); //writing the data to the memory buffer
65
66
67

```

```

66     comp_buff_position += sprintf(comp_output + comp_buff_position, "%zu%c", counter, character); //writing the data to the memory buffer
67
68
69
70     i += counter; // Adjusting i variable for the loop to progress through the characters
71 }
72
73
74 seg_data->output_buffer = comp_output; //Writing the output data to output buffer
75 seg_data->output_length = comp_buff_position; //Adjusting the position of buffer
76
77
78
79 pthread_exit(NULL); //exiting thread
80
81 }
82
83
84
85
86
87
88
89
90
91 //-----Main function-----
92
93
94
95
96 int main(int arg_counter, char *arg_select[]) { //establishing the main function
97
98
99
100     if (arg_counter < 2) { //Handling incorrect usage cases by throwing usage suggestion
101
102         message_printer("usage: pzip <file1> <file2> ... > <outputFile>"); //calling message printer to print correct usage
103
104
105         exit(1); //exiting the program
106     }
107
108
109
110
111     //Next step is to establish the threads, doing it based on the number of processors system has
112
113     int th_number = sysconf(_SC_NPROCESSORS_ONLN); //getting number of threads by getting the number of processors of the system. Using sysconf f
114
115
116
117     if (th_number < 1) { //in case sysconf fails, we set the number of threads to 1 as default number (1 thread)
118         th_number = 1;
119     }
120
121
122     pthread_t threads[th_number]; //Initializing the threads based on the established amount
123
124     Th_data th_data[th_number]; //Establishing Th_data instances (amount of instances based on thread number). Th_data data structure has been in
125
126
127
128
129     for (int i = 1; i < arg_counter; i++) { //establishing a loop for all the arguments (input files)
130
131         int input_file = open(arg_select[i], O_RDONLY); //opening the file (one by one as it's a loop)

```

```
for (int i = 1; i < arg_counter; i++) { //establishing a loop for all the arguments (input files)
```

```
    int input_file = open(arg_select[i], O_RDONLY); //opening the file (one by one as it's a loop)
```

```
    if (input_file < 0) { //Error handling for file opening
```

```
        message_printer("Could not open the file");
```

```
        exit(1);
```

```
    }
```

```
    struct stat file_stats; //initializing instance for collecting file statistics (for example size)
```

```
    if (fstat(input_file, &file_stats) != 0) { //Error handling for file statistics
```

```
        message_printer("Couldn't get the size of the file");
```

```
        close(input_file);
```

```
        exit(1);
```

```
    }
```

```
    char *data = mmap(NULL, file_stats.st_size, PROT_READ, MAP_PRIVATE, input_file, 0); //mapping file data into memory for easy access
```

```
    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 each
```

```
    for (int k = 0; k < th_number; k++) { //Using for-loop to create each thread 1 by 1, for file compression
```

```
        //Defining each element of the thread's data structure, such as start of segment, segment length and thread index
```

```
        th_data[k].seg_start = data + k * seg_size; //Defining start of segment. The calculation is an offset that the thread has to make to get
```

```
        th_data[k].seg_length = (k == th_number - 1) ? (file_stats.st_size - k * seg_size) : seg_size; //Defining length of segment, includes last
```

```
        th_data[k].th_index = k; //setting index of thread as k number of the for-loop (simple counting)
```

```
        pthread_create(&threads[k], NULL, seg_compression, &th_data[k]); //creating thread, function for compressing the segment specified
```

```
    }
```

```
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 the file

    free(th_data[j].output_buffer); //Freeing the output buffer
}

munmap(data, file_stats.st_size); //Freeing mapped memory
close(input_file); //closing input file
}

return 0;
}
```


Screenshots of the punzip.c code:

```
C punzip.c > Th_data
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;
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
36
37
38
39
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
48          message_printer("malloc failed");
49
50          pthread_exit(NULL);
51      }
52
53      size_t decomp_buff_position = 0; //initializing variable for positioning in the buffer
54
55      for (size_t i = 0; i < seg_length; i++) { //Using for-loop to go over each encoded character inside the segment
56
57          int c_count; //Initializing count value for RLE encoding.
58          char enc_character; //Initializing character value for RLE encoding.
59
60
61          if (sscanf(seg_start + i, "%d%c", &c_count, &enc_character) == 2) { //checking if the encoding format is correct
62
63              //printf("%d %zu\n", c_count, th_index); //print statement for troubleshooting, can safely be removed
64
65              for (int j = 0; j < c_count; j++) { //Using for-loop to get the decoded output
66
67
68                  decomp_output[decomp_buff_position++] = enc_character;
```

```

64
65     for (int j = 0; j < c_count; j++) {          //Using for-loop to get the decoded output
66
67
68         decomp_output[decomp_buff_position++] = enc_character;
69
70     }
71
72     while (i < seg_length && isdigit(seg_start[i])) i++; //Skip the count of the pair
73     i++; // Skip the character from the pair
74
75 } else { //error handling for format errors. Throws error if program detects incorrect encoding format
76
77     message_printer("Wrong encoding format"); //throwing the error
78
79     pthread_exit(NULL); //exiting the thread
80
81 }
82
83 }
84
85
86 seg_data->output_buffer = decomp_output; //Writing the result into the buffer
87 seg_data->output_length = decomp_buff_position; //Updating the buffer position
88
89
90
91 pthread_exit(NULL); //exiting thread
92
93 }
94
95
96
97
98
99
100
101
102
103 //-----Main function-----
104
105
106
107
108
109
110 int main(int arg_counter, char *arg_select[]) { //establishing the main function
111
112
113
114     if (arg_counter < 2) {          //Handling incorrect usage cases by throwing usage suggestion
115
116         message_printer("usage: punzip <file1> <file2> ... > <outputFile1> <outputFile2>"); //calling message printer to print correct usage
117
118
119         exit(1); //exiting the program
120     }
121
122
123
124
125 //Next step is to establish the threads, doing it based on the number of processors system has
126
127 int th_number = sysconf(_SC_NPROCESSORS_ONLN); //getting number of threads by getting the number of processors of the system. Using sysconf fun
128
129
130
131 if (th_number < 1) { //in case sysconf fails, we set the number of threads to 1 as default number (1 thread)
132     th_number = 1;
133 }
134
135
136 pthread_t threads[th_number]; //Initializing the threads based on the established amount

```

```

134
135
136 pthread_t threads[th_number]; //Initializing the threads based on the established amount
137
138 Th_data th_data[th_number]; //Establishing Th_data instances (amount of instances based on thread number). Th_data data structure has been init
139
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
160
161     if (fstat(input_file, &file_stats) != 0) { //Error handling for file statistics
162
163         message_printer("Couldn't get the size of the file");
164
165         close(input_file);
166
167         exit(1);
168     }
169
170
171
172     char *data = mmap(NULL, file_stats.st_size, PROT_READ, MAP_PRIVATE, input_file, 0); //mapping file data into memory for easy access
173
174     if (data == MAP_FAILED) { //Error handling for file data mapping
175
176         message_printer("Couldn't map the file data to memory");
177
178         close(input_file);
179
180         exit(1);
181     }
182
183
184
185     size_t seg_size = file_stats.st_size / th_number; //establishing segments' sizes based on file size divided by number of threads, so that e
186
187
188
189     for (int k = 0; k < th_number; k++) {
190
191         size_t seg_start_offset = k * seg_size; //Getting start offset for the thread to know where to start
192         size_t seg_end = (k == th_number - 1) ? file_stats.st_size : seg_start_offset + seg_size; //Determining the end of the segment
193
194
195
196         if (k > 0) { //Using if-conditions and a loop to adjust start offset for the segments that are following after the 1st one
197
198             while (seg_start_offset < seg_end && isdigit(data[seg_start_offset])) {
199                 seg_start_offset++;
200             }
201
202             if (seg_start_offset < seg_end) { //Using if condition for the offset to move to a full pair or count-character

```

```

202     if (seg_start_offset < seg_end) { //Using if condition for the offset to move to a full pair or count-character
203         seg_start_offset++;
204     }
205 }
206
207
208 while (seg_end < file_stats.st_size && isdigit(data[seg_end])) { //Using a loop and if condition to adjust the end of the segment as we
209     seg_end++;
210 }
211 if (seg_end < file_stats.st_size) {
212     seg_end++;
213 }
214
215
216
217 if (seg_start_offset >= file_stats.st_size || seg_end > file_stats.st_size) { //Adding error handling for segment boundaries
218     message_printer("Out of bounds error in segment. The compressed file may be corrupt");
219     exit(1);
220 }
221
222 //Establishing thread data for creating the thread
223 th_data[k].seg_start = data + seg_start_offset;
224 th_data[k].seg_length = seg_end - seg_start_offset;
225 th_data[k].th_index = k;
226
227 //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);
228
229 pthread_create(&threads[k], NULL, seg_decompression, &th_data[k]); //Creating thread
230 }
231
232

```

```

233
234
235 for (int j = 0; j < th_number; j++) { //Using for-loop for joining all the threads
236
237     pthread_join(threads[j], NULL);
238
239     fwrite(th_data[j].output_buffer, sizeof(char), th_data[j].output_length, stdout); //Writing to a file
240
241     free(th_data[j].output_buffer); //Freeing the output buffer
242 }
243
244
245
246
247 munmap(data, file_stats.st_size); //Freeing mapped memory
248
249 close(input_file); //closing input file
250 }
251
252
253
254
255
256
257 return 0;
258
259 }
260
261

```