

# Теория кодирования и сжатия информации

## Лабораторная работа №6

Гущин Андрей, 431 группа, 1 подгруппа

2022 г.

### 1 Задача

Разработать программу осуществляющую архивацию и разархивацию текстового файла используя алгоритм Гилберта-Мура. Программы архивации и разархивации должны быть представлены отдельно и работать независимо друг от друга. Определить для данного шифра характеристики 1 (коэффициент сжатия) и 2 (скорость сжатия). К работе необходимо прикрепить отчет и программный проект.

### 2 Алгоритм

Алгоритм Гилберта-Мура заключается в создании нового кода для каждого встречающегося в тексте символа на основе частоты встречи этого символа.

Алгоритм состоит из следующих шагов:

1. Вычислить частоты  $p_i$  всех встретившихся символов;
2. Отсортировать полученный список по невозрастанию;
3. Вычислить префиксные суммы  $q_i$  для списка частот;
4. Вычислить  $\sigma_i = q_i + p_i/2$
5. Перевести  $\sigma_i$  из десятичной системы счисления в двоичную;
6. В качестве кода для символа  $i$  необходимо использовать первые  $\lceil -\log_2 p_i \rceil + 1$  двоичных цифр после запятой из  $\sigma_i$ .

### 3 Тестирование

Для проверки программы были использованы тестовые тексты 1 (рис. 1) и 6 (рис. 2). Можно заметить, что после распаковки архива полученный файл совпадает с исходным (проверка с помощью утилиты diff). Также можно заметить, что для файлов малого размера архив увеличивает их размер за счёт метаданных.

```
-zsh
$ ./lab6
error: The following required arguments were not provided:
  -i <INPUT_FILE>
  -o <OUTPUT_FILE>

Usage: lab6 -i <INPUT_FILE> -o <OUTPUT_FILE>

For more information try '--help'
$ ./lab6 -i test1.txt -o test1.gs6 --compress
$ ./lab6 -i test1.gs6 -o test1.out.txt --decompress
$ ls -l
total 1832
-rwxr-xr-x  1 aguschin  staff   924546 Dec 15 13:47 lab6
-rw-r--r--  1 aguschin  staff      13 Dec 15 13:51 test1.gs6
-rw-r--r--  1 aguschin  staff        2 Dec 15 13:51 test1.out.txt
-rw-----@ 1 aguschin  staff        2 Dec 15 13:47 test1.txt
$ diff --report-identical-files test1.out.txt test1.txt
Files test1.out.txt and test1.txt are identical
$
```

Рис. 1: Сжатие текста Тест\_1.txt

```
-zsh
$ ./lab6
error: The following required arguments were not provided:
  -i <INPUT_FILE>
  -o <OUTPUT_FILE>

Usage: lab6 -i <INPUT_FILE> -o <OUTPUT_FILE>

For more information try '--help'
$ ./lab6 -i test6.txt -o test6.gs6 --compress
$ ./lab6 -i test6.gs6 -o test6.out.txt --decompress
$ ls -l
total 1856
-rwxr-xr-x  1 aguschin  staff   924546 Dec 15 13:47 lab6
-rw-r--r--  1 aguschin  staff    7337 Dec 15 13:53 test6.gs6
-rw-r--r--  1 aguschin  staff    7958 Dec 15 13:54 test6.out.txt
-rw-----@ 1 aguschin  staff    7958 Dec 15 13:53 test6.txt
$ diff --report-identical-files test6.txt test6.out.txt
Files test6.txt and test6.out.txt are identical
$
```

Рис. 2: Сжатие текста Тест\_6.txt

## 4 Вычисленные характеристики

### 4.1 Характеристика 1 (Коэффициент сжатия)

Результаты применения программы к каждому из тестовых текстовых файлов занесены в таблицу 1.

Название	Исходный размер, байт	Сжатый размер, байт	Коэффициент
Тест_1.txt	2	13	0.15385
Тест_2.txt	33	196	0.16837
Тест_3.txt	2739	2564	1.06825
Тест_4.txt	330	95	3.47368
Тест_5.txt	59	349	0.16905
Тест_6.txt	7958	7337	1.08464
Тест_7.txt	138245	108673	1.27212
Тест_8.txt	574426	435482	1.31906
Тест_9.txt	2752	352	7.81818
Тест_10.txt	2814	768	3.66406

Таблица 1: результаты тестирования

### 4.2 Характеристика 2 (Скорость сжатия)

Для тестирования скорости сжатия использовался произвольный двоичный файл размера 76450438 байт ( $\approx 73$  мегабайта). В результате пяти последовательных запусков, среднее время запаковки файла составило 3.68 секунд, среднее время распаковки составило 180 секунд.

Таким образом, средняя скорость сжатия составила 19.81218 Мбайт в секунду, а средняя скорость разжатия составила 0.40505 Мбайт в секунду.

## 5 Реализация

Программа реализована на языке программирования Rust с использованием библиотеки clap для чтения параметров командной строки. Сборка производится с помощью программы cargo, поставляющейся вместе с языком.

### 5.1 Содержимое файла gilbert\_moore.rs

```
1 use std::cmp::Ordering;
2 use std::collections::HashMap;
3
4 #[derive(Debug, Clone)]
5 pub struct Weighted<T> {
6     pub weight: f32,
7     pub value: T,
8 }
9
10 #[derive(Clone, Copy, Debug, PartialEq)]
11 enum Bit {
12     Zero,
```

```

13     One,
14 }
15 impl Bit {
16     fn from_u8(x: u8) -> Bit {
17         if x == 0 {
18             Bit::Zero
19         } else {
20             Bit::One
21         }
22     }
23
24     fn to_u8(x: &Bit) -> u8 {
25         match *x {
26             Bit::Zero => 0,
27             Bit::One  => 1,
28         }
29     }
30 }
31
32 fn get_probabilities(data: &Vec<u8>) -> Vec<Weighted<u8>> {
33     let mut counts = HashMap::new();
34     data.iter().for_each(|byte| {
35         *counts.entry(*byte).or_insert(0) += 1;
36     });
37     let mut counts: Vec<(u8, u32)> = counts.drain().collect();
38     counts.sort();
39
40     let total = data.len() as f32;
41     let mut weights = Vec::new();
42     counts.iter().for_each(|(byte, count)| {
43         weights.push(Weighted {
44             value: *byte,
45             weight: *count as f32 / total,
46         });
47     });
48     return weights;
49 }
50
51 fn prefix_sum(weights: &Vec<Weighted<u8>>) -> Vec<f32> {
52     let mut pf = Vec::new();
53     for i in 0..weights.len() {
54         if i == 0 {
55             pf.push(0f32);
56         } else {
57             pf.push(weights[i - 1].weight + pf[i - 1]);
58         }
59     }
60     return pf;
61 }
62
63 fn get_code(probabilities: &Vec<Weighted<u8>>) -> HashMap<u8, Vec<Bit>> {
64     let mut code = HashMap::new();
65
66     let mut ps = probabilities.clone();

```

```

67     ps.sort_by(|a, b|
↪     a.weight.partial_cmp(&b.weight).unwrap_or(Ordering::Equal));
68     ps.reverse();
69     let prefix = prefix_sum(&ps);
70
71     for i in 0..ps.len() {
72         let l = (-ps[i].weight.log2()).ceil() as i32 + 1;
73         if l == 1 {
74             code.insert(ps[i].value, vec![Bit::Zero]);
75         } else {
76             let sigma = prefix[i] + ps[i].weight as f32 / 2f32;
77             let bits = get_bits(sigma, l);
78             code.insert(ps[i].value, bits);
79         }
80     }
81
82     return code;
83 }
84
85 fn get_bits(x: f32, count: i32) -> Vec<Bit> {
86     let mut bits = Vec::new();
87
88     let mut tmp = x;
89     for _ in 0..count {
90         tmp *= 2f32;
91         if tmp >= 1f32 {
92             bits.push(Bit::One);
93             tmp -= 1f32;
94         } else {
95             bits.push(Bit::Zero);
96         }
97     }
98
99     return bits;
100 }
101
102 #[derive(Debug)]
103 struct Metadata {
104     probabilities: Vec<Weighted<u8>>,
105     code: HashMap<u8, Vec<Bit>>,
106     remainder: u8,
107 }
108
109 impl Metadata {
110     fn compute(data: &Vec<u8>) -> Self {
111         let probs = get_probabilities(&data);
112         let code = get_code(&probs);
113
114         return Self {
115             probabilities: probs,
116             code,
117             remainder: 0,
118         };
119     }

```

```

120
121 fn load(data: &Vec<u8>) -> Self {
122     let remainder = data[0];
123     let dict_len = data[1] as usize + 1;
124     let mut probabilities = Vec::new();
125     for i in 0..dict_len {
126         let pstart = 2 + 5 * i + 1;
127         let prob = [
128             data[pstart + 0],
129             data[pstart + 1],
130             data[pstart + 2],
131             data[pstart + 3],
132         ];
133         let prob: f32 = unsafe { std::mem::transmute(prob) };
134         probabilities.push(Weighted {
135             value: data[2 + 5 * i],
136             weight: prob,
137         });
138     }
139     let code = get_code(&probabilities);
140
141     return Self {
142         probabilities,
143         code,
144         remainder,
145     };
146 }
147
148 fn dump(&self) -> Vec<u8> {
149     let mut result = Vec::new();
150
151     result.push(self.remainder);
152     result.push((self.probabilities.len() - 1) as u8);
153     for p in &self.probabilities {
154         result.push(p.value);
155         let bweight: [u8; 4] = unsafe { std::mem::transmute(p.weight)
↪ };
156         for b in bweight {
157             result.push(b);
158         }
159     }
160
161     return result;
162 }
163 }
164
165 struct BitWriter {
166     buffer: Vec<Bit>,
167     remainder: u8,
168     result: Vec<u8>,
169 }
170
171 impl BitWriter {
172     fn new() -> Self {

```

```

173         Self {
174             buffer: Vec::new(),
175             remainder: 0,
176             result: Vec::new(),
177         }
178     }
179
180     fn dump_byte(&mut self) {
181         let mut byte = 0_u8;
182         self.buffer
183             .iter()
184             .map(Bit::to_u8)
185             .enumerate()
186             .for_each(|(i, bit)| byte |= bit << i);
187         self.result.push(byte);
188         self.buffer.clear();
189     }
190
191     fn write_bit(&mut self, bit: Bit) {
192         self.buffer.push(bit);
193         if self.buffer.len() == 8 {
194             self.dump_byte();
195         }
196     }
197
198     fn write_bits(&mut self, bits: &Vec<Bit>) {
199         bits.iter().for_each(|bit| self.write_bit(*bit));
200     }
201
202     fn finish(&mut self) {
203         let remainder = 8 - self.buffer.len() as u8 % 8;
204         if remainder != 0 {
205             self.dump_byte();
206         }
207         self.remainder = remainder;
208     }
209 }
210
211 pub fn compress(data: &Vec<u8>) -> Vec<u8> {
212     let mut result = Vec::new();
213     let mut metadata = Metadata::compute(data);
214
215     let mut writer = BitWriter::new();
216     data.iter().for_each(|byte| {
217         let bits = metadata.code.get(byte).unwrap();
218         writer.write_bits(&bits);
219     });
220     writer.finish();
221     metadata.remainder = writer.remainder;
222
223     let md_dump = metadata.dump();
224     md_dump.iter().for_each(|byte| result.push(*byte));
225     writer.result.iter().for_each(|byte| result.push(*byte));
226

```

```

227     return result;
228 }
229
230 struct BitReader<'a> {
231     data: &'a Vec<u8>,
232     metadata: &'a Metadata,
233     buffer: Vec<Bit>,
234     ptr: usize,
235 }
236
237 impl<'a> BitReader<'a> {
238     fn new(data: &'a Vec<u8>, metadata: &'a Metadata) -> Self {
239         Self {
240             data,
241             metadata,
242             buffer: Vec::new(),
243             ptr: 0,
244         }
245     }
246
247     fn read_byte(&mut self) {
248         if let Some(byte) = self.data.get(self.ptr) {
249             for i in 0..=7 {
250                 let bit = byte & (1 << i);
251                 self.buffer.push(Bit::from_u8(bit));
252             }
253             if self.ptr == self.data.len() - 1 {
254                 for _ in 0..self.metadata.remainder {
255                     self.buffer.pop();
256                 }
257             }
258             self.buffer.reverse();
259             self.ptr += 1;
260         }
261     }
262
263     fn read_bit(&mut self) -> Option<Bit> {
264         if self.buffer.len() == 0 {
265             self.read_byte();
266         }
267         self.buffer.pop()
268     }
269 }
270
271 pub fn decompress(archive: &Vec<u8>) -> Vec<u8> {
272     let mut result = Vec::new();
273     let metadata = Metadata::load(archive);
274
275     let data = archive[2 + metadata.probabilities.len() * 5..].to_vec();
276     let mut reader = BitReader::new(&data, &metadata);
277
278     let mut run = Vec::new();
279     while let Some(bit) = reader.read_bit() {
280         run.push(bit);

```



```

281         'checker: for p in &metadata.code {
282             if p.1.len() != run.len() {
283                 continue;
284             }
285             for i in 0..run.len() {
286                 if p.1[i] != run[i] {
287                     continue 'checker;
288                 }
289             }
290             result.push(*p.0);
291             run.clear();
292             break;
293         }
294     }
295
296     return result;
297 }

```

## 5.2 Содержимое файла main.rs

```

1  mod gilbert_moore;
2  use clap::Parser;
3  use std::fs::File;
4  use std::io::{Error, ErrorKind, Read, Write};
5  use std::path::PathBuf;
6
7  #[derive(Parser)]
8  struct Cli {
9      #[arg(short)]
10     input_file: PathBuf,
11
12     #[arg(short)]
13     output_file: PathBuf,
14
15     #[arg(long, default_value_t = true)]
16     compress: bool,
17
18     #[arg(long, default_value_t = false)]
19     decompress: bool,
20 }
21
22 fn run_decompressor(cli: &Cli) -> Result<(), Error> {
23     let mut input_f = File::open(cli.input_file.to_str().unwrap())?;
24     let mut archive = Vec::new();
25     input_f.read_to_end(&mut archive)?;
26     let data = gilbert_moore::decompress(&archive);
27
28     let mut output_f = File::create(cli.output_file.to_str().unwrap())?;
29     output_f.write_all(&data)?;
30
31     return Ok(());
32 }
33
34 fn run_compressor(cli: &Cli) -> Result<(), Error> {

```

```

35     let mut input_f = File::open(cli.input_file.to_str().unwrap())?;
36     let mut data = Vec::new();
37     input_f.read_to_end(&mut data)?;
38     let archive = gilbert_moore::compress(&data);
39
40     let mut output_f = File::create(cli.output_file.to_str().unwrap())?;
41     output_f.write_all(&archive)?;
42
43     return Ok(());
44 }
45
46 fn main() -> std::io::Result<()> {
47     let cli = Cli::parse();
48
49     let result = if cli.decompress {
50         run_decompressor(&cli)
51     } else {
52         run_compressor(&cli)
53     };
54
55     if let Err(error) = result {
56         match error.kind() {
57             ErrorKind::NotFound => println!("Указанный файл не найден"),
58             ErrorKind::AlreadyExists => println!("Указанный файл уже
↪ существует"),
59             _ => println!("Произошла непредвиденная ошибка"),
60         };
61     }
62
63     Ok(())
64 }

```

### 5.3 Содержимое файла Cargo.toml

```

1  [package]
2  name = "lab6"
3  version = "0.1.0"
4  edition = "2021"
5
6  [dependencies]
7  clap = { version = "4.0.17", features = ["derive"] }

```