

Теория кодирования и сжатия информации

Лабораторная работа №8

Гущин Андрей, 431 группа, 1 подгруппа

2022 г.

1 Задача

Разработать программу осуществляющую архивацию и разархивацию текстового файла используя алгоритм Лемпеля-Зива (LZ78). 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 LZ78 заключается в создании нового кода для всех встречающихся последовательностей символов. Алгоритм использует динамический словарь.

Алгоритм состоит из следующих шагов:

1. Прочитать первый символ. Так как он ни разу не встречался, добавить его в словарь.
2. Далее необходимо найти наибольшую подстроку ps , уже находящуюся в словаре.
3. Если строка не найдена, то просто добавить отдельный символ аналогично самому первому. Иначе записать найденную подстроку и добавить в словарь конкатенацию найденной записи и следующего символа.

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.

```

$ ./lab8
error: The following required arguments were not provided:
  -i <INPUT_FILE>
  -o <OUTPUT_FILE>

Usage: lab8 -i <INPUT_FILE> -o <OUTPUT_FILE>

For more information try '--help'
$ ./lab8 -i test1.txt -o test1.gs8 --compress
$ ./lab8 -i test1.gs8 -o test1.out.txt --decompress
$ ls -l
total 1792
-rwxr-xr-x  1 aguschin  staff   904306 Dec 15 14:41 lab8
-rw-r--r--  1 aguschin  staff      6 Dec 15 14:41 test1.gs8
-rw-r--r--  1 aguschin  staff      2 Dec 15 14:42 test1.out.txt
-rw-----@ 1 aguschin  staff      2 Dec 15 14:41 test1.txt
$ diff --report-identical-files test1.txt test1.out.txt
Files test1.txt and test1.out.txt are identical
$
```

Рис. 1: Сжатие текста Тест_1.txt

```

$ ./lab8
error: The following required arguments were not provided:
  -i <INPUT_FILE>
  -o <OUTPUT_FILE>

Usage: lab8 -i <INPUT_FILE> -o <OUTPUT_FILE>

For more information try '--help'
$ ./lab8 -i test6.txt -o test6.gs8 --compress
$ ./lab8 -i test6.gs8 -o test6.out.txt --decompress
$ ls -l
total 1816
-rwxr-xr-x  1 aguschin  staff   904306 Dec 15 14:42 lab8
-rw-r--r--  1 aguschin  staff    7683 Dec 15 14:42 test6.gs8
-rw-r--r--  1 aguschin  staff   7958 Dec 15 14:43 test6.out.txt
-rw-----@ 1 aguschin  staff   7958 Dec 15 14:42 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	6	0.33333
Тест_2.txt	33	99	0.33333
Тест_3.txt	2739	1227	2.23227
Тест_4.txt	330	108	3.05556
Тест_5.txt	59	177	0.33333
Тест_6.txt	7958	7683	1.03579
Тест_7.txt	138245	105015	1.31643
Тест_8.txt	574426	405459	1.41673
Тест_9.txt	2752	222	12.3964
Тест_10.txt	2814	351	8.01709

Таблица 1: результаты тестирования

4.2 Характеристика 2 (Скорость сжатия)

Для тестирования скорости сжатия использовался произвольный двоичный файл размера 3120002 байт (≈ 3 мегабайта). В результате пяти последовательных запусков, среднее время запаковки файла составило 11.28 секунды, среднее время распаковки составило 0.02 секунд.

Таким образом, средняя скорость сжатия составила 0.26378 Мбайт в секунду, а средняя скорость распаковки составила 148.77329 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 lz78.rs

```
1  const DICT_SIZE: usize = 4096;
2
3  pub fn lz_match(data: &Vec<u8>, pos1: usize, pos2: usize, length: usize)
4  ↪ -> bool {
5      for i in 0..length {
6          if data[pos1 + i] != data[pos2 + i] {
7              return false;
8          }
9      }
10     return true;
11 }
```

```

12 fn lz78_encode(data: &Vec<u8>) -> Vec<(u32, u8)> {
13     let mut encoded = Vec::new();
14     let mut dict = Vec::new();
15     let mut count = 1;
16
17     let mut ptr = 0;
18     while ptr < data.len() {
19         let mut saved = None;
20         let mut l = 1;
21         while ptr + l < data.len() && l < 256 {
22             let mut found = false;
23             for i in (0..dict.len()).rev() {
24                 let (pos1, l1, val) = dict[i];
25                 if l != l1 {
26                     continue;
27                 }
28                 if lz_match(data, pos1, ptr, l) {
29                     found = true;
30                     saved = Some((pos1, l, val));
31                     l += 1;
32                     break;
33                 }
34             }
35             if !found {
36                 break;
37             }
38         }
39         if let Some((_, l, val)) = saved {
40             let t = (val, data[ptr + l]);
41             encoded.push(t);
42             if dict.len() < DICT_SIZE {
43                 dict.push((ptr, l + 1, count));
44                 count += 1;
45             }
46             ptr += l + 1;
47         } else {
48             encoded.push((0, data[ptr]));
49             if dict.len() < DICT_SIZE {
50                 dict.push((ptr, 1, count));
51                 count += 1;
52             }
53             ptr += 1;
54         }
55     }
56
57     return encoded;
58 }
59
60 pub fn compress(data: &Vec<u8>) -> Vec<u8> {
61     let mut result = Vec::new();
62     let encoded = lz78_encode(data);
63
64     for (value, byte) in &encoded {

```

```

65         let value_bytes: [u8; 4] = unsafe { std::mem::transmute(*value)
↳ };
66         result.push(value_bytes[0]);
67         result.push(value_bytes[1]);
68         result.push(*byte);
69     }
70
71     return result;
72 }
73
74 pub fn decompress(data: &Vec<u8>) -> Vec<u8> {
75     let mut decoded = Vec::new();
76     let mut dict = Vec::new();
77     let mut count = 1;
78     let mut ptr = 0;
79
80     let mut caret = 0;
81     while caret < data.len() {
82         let value_bytes: [u8; 4] = [data[caret + 0], data[caret + 1], 0,
↳ 0];
83         let value: u32 = unsafe { std::mem::transmute(value_bytes) };
84         let byte = data[caret + 2];
85         caret += 3;
86
87         if value == 0 {
88             if dict.len() < DICT_SIZE {
89                 dict.push((ptr, 1, count));
90                 count += 1;
91             }
92             decoded.push(byte);
93             ptr += 1;
94         } else {
95             let (pos1, l, _) = dict[value as usize - 1];
96             for i in 0..l {
97                 decoded.push(decoded[pos1 + i]);
98             }
99             decoded.push(byte);
100             if dict.len() < DICT_SIZE {
101                 dict.push((ptr, l + 1, count));
102                 count += 1;
103             }
104             ptr += l + 1;
105         }
106     }
107
108     return decoded;
109 }

```

5.2 Содержимое файла main.rs

```

1 mod lz78;
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 = lz78::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 = lz78::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 = "lab8"
3  version = "0.1.0"
4  edition = "2021"
5
6  [dependencies]
7  clap = { version = "4.0.17", features = ["derive"] }

```