

НАЦІОНАЛЬНИЙ ТЕХНІЧНИЙ УНІВЕРСИТЕТ УКРАЇНИ
«КІЇВСЬКИЙ ПОЛІТЕХНІЧНИЙ ІНСТИТУТ ім. Ігоря СІКОРСЬКОГО»
ФІЗИКО-ТЕХНІЧНИЙ ІНСТИТУТ

Комп'ютерний практикум №4
з курсу алгоритми кодування двійкових
даних

**РЕАЛІЗАЦІЯ АЛГОРИТМУ
СТИСКАННЯ ДАНИХ LZW**

Виконав студент
групи ФІ-42МН
Бєш Радомир Андрійович

Зміст

1	Мета	3
2	Постановка задачі	3
3	Хід роботи	3
3.1	Особливості реалізації алгоритму LZW	3
3.2	Порівняльний аналіз	6
4	Висновки	8

1 Мета

Опанувати методи роботи із двійковими потоками даних та алгоритми кодування даних зі змінними та/або нефіксованими довжинами кодових слів.

2 Постановка задачі

У даній роботі необхідно реалізувати алгоритм стискання даних Зіва-Лемпеля-Велч (LZW). Реалізація повинна бути у вигляді автономної утиліти, яка дозволяє стискати файли у архіви та розпаковувати їх назад. Для полегшення вважаємо, що реалізація працює із файлами, розмір яких не перевищує 2^{32} байтів (приблизно 4 Гб).

Результати дослідження:

https://github.com/Radomir21/Encoding-Algorithms/tree/main/lab_4.

3 Хід роботи

3.1 Особливості реалізації алгоритму LZW

У моїй лабораторній роботі використовується модуль BitStream побітового читання і запису даних з файлу, написаний в ЛБ №2.

1. Організація словника алгоритму LZW:

Алгоритм використовує словник, який реалізовано у вигляді звичайного словника Python (хеш-таблиці). Початковий словник містить усі можливі одиничні байти (0–255), яким відповідають початкові коди. Пошук фраз у словнику виконується шляхом прямої перевірки наявності ключа у словнику. Якщо така фраза вже є у словнику, вона розширюється. Якщо ж ні, то код поточної фрази записується у бітовий потік, а нова фраза додається до словника з наступним вільним кодом.

```

53     def lzw_compress(input_path, output_path=None, max_bits=16, overflow_mode="reset"):
54         max_dict_size = 2^max_bits
55
56         if output_path is None:
57             output_path = input_path + ".lzw"
58
59         with open(input_path, "rb") as fin, open(output_path, "wb") as fout:
60             write_header(fout, max_bits, overflow_mode)
61             bs = BitStream(fout, "w")
62
63             dictionary = {bytes([i]): i for i in range(256)}
64             next_code = 256
65
66             w = b""
67
68             while True:
69                 chunk = fin.read(4096)
70                 if not chunk:
71                     break
72
73                 for byte_val in chunk:
74                     k = bytes([byte_val])
75                     wk = w + k
76
77                     if wk in dictionary:
78                         w = wk
79                     else:
80                         if w:
81                             write_bits_fixed(bs, dictionary[w], max_bits)
82
83                         if next_code < max_dict_size:
84                             dictionary[wk] = next_code
85                             next_code += 1
86                         else:
87                             if overflow_mode == "reset":
88                                 dictionary = {bytes([i]): i for i in range(256)}
89                                 next_code = 256
90
91                         w = k
92
93                     if w:
94                         write_bits_fixed(bs, dictionary[w], max_bits)
95
96             bs.close()
97
98         return output_path

```

2. Алгоритм стиснення LZW:

Стиснення виконує функція `lzw_compress()`. Алгоритм побайтово зчитує вхідний файл та на кожному кроці намагається знайти у словнику найдовшу фразу, що вже там присутня. Для цього формується нова послідовність $w + k$, де w — поточна фраза, а k — наступний байт.

```

73             for byte_val in chunk:
74                 k = bytes([byte_val])
75                 wk = w + k
76
77                 if wk in dictionary:
78                     w = wk
79                 else:
80                     if w:
81                         write_bits_fixed(bs, dictionary[w], max_bits)

```

3. Обробка переповнення словника:

У реалізації передбачено два режими обробки переповнення словника: reset — словник очищується та повертається до початкового стану, freeze — словник перестає розширюватися. Вибір режиму здійснюється за допомогою параметра та зберігається у хедері стисненого файлу для коректного декодування.

```

83     if next_code < max_dict_size:
84         dictionary[wk] = next_code
85         next_code += 1
86     else:
87         if overflow_mode == "reset":
88             dictionary = {bytes([i]): i for i in range(256)}
89             next_code = 256
90

```

4. Заголовок стисненого файлу:

Для коректного декодування у стисненому файлі зберігається заголовок, який містить сигнатуру формату (MAGIC) - це просто фіксована послідовність байтів, записана на початку файлу, максимальну довжину кодового слова та режим переповнення словника. Заголовок записується перед стисненими даними та читається під час декодування. Заголовок — це кілька байтів на початку файлу, які описують формат файлу та передають параметри кодування декодеру.

```

30     def write_header(out_f, max_bits, overflow_mode):
31         out_f.write(MAGIC)
32         write_uint8(out_f, max_bits)
33
34         mode_byte = 0 if overflow_mode == "freeze" else 1
35         write_uint8(out_f, mode_byte)
36
37
38     def read_header(in_f):
39         magic = in_f.read(4)
40         if magic != MAGIC:
41             raise ValueError("Це не LZW файл (невірний MAGIC)")
42
43         max_bits = read_uint8(in_f)
44         mode_byte = read_uint8(in_f)
45
46         if max_bits < 9 or max_bits > 16:
47             raise ValueError(f"Некоректне max_bits={max_bits} (очікується 9..16)")
48
49         overflow_mode = "freeze" if mode_byte == 0 else "reset"
50         return max_bits, overflow_mode
51

```

5. Алгоритм декодування:

Декодування реалізовано у функції `lzw_decompress`. Декодер полягає у послідовному зчитуванні кодів зі стисненого файлу та заміні їх відповідними послідовностями байтів зі словника. Словник формується динамічно під час декодування у тому ж порядку, що і при стисненні. У випадку, коли зчитаний код ще відсутній у словнику, фраза відновлюється як конкатенація попередньої фрази та її першого байта.

```

101     def lzw_decompress(input_path, output_path=None):
102         if output_path is None:
103             if input_path.endswith(".lzw"):
104                 output_path = input_path[:-4]
105             else:
106                 output_path = input_path + ".out"
107
108         with open(input_path, "rb") as fin, open(output_path, "wb") as fout:
109             max_bits, overflow_mode = read_header(fin)
110             max_dict_size = 1 << max_bits
111
112             bs = BitStream(fin, "r")
113
114             dictionary = {i: bytes([i]) for i in range(256)}
115             next_code = 256
116
117             try:
118                 prev_code = read_bits_fixed(bs, max_bits)
119             except EOFError:
120                 return output_path
121
122             if prev_code not in dictionary:
123                 raise ValueError("Пошкоджені дані: перший код не в словнику")
124
125             w = dictionary[prev_code]
126             fout.write(w)
127
128             while True:
129                 try:
130                     code = read_bits_fixed(bs, max_bits)
131                 except EOFError:
132                     break
133
134                 if code in dictionary:
135                     entry = dictionary[code]
136                 elif code == next_code:
137                     entry = w + w[:1]
138                 else:
139                     raise ValueError(f"Пошкоджені дані: зустріли невірний код {code}")
140
141                 fout.write(entry)
142
143                 if next_code < max_dict_size:
144                     dictionary[next_code] = w + entry[:1]
145                     next_code += 1
146                 else:
147                     if overflow_mode == "reset":
148                         dictionary = {i: bytes([i]) for i in range(256)}
149                         next_code = 256
150
151                     w = entry
152
153             bs.close()
154
155         return output_path
156

```

3.2 Порівняльний аналіз

Для аналізу було використано такі типи файлів по 10 файлів кожного типу:

- **TXT**
- **CSV**
- **BMP**
- **PDF**
- **EXE**

Коефіцієнт стискання розраховується, як:

$$K = \frac{\text{SIZE}_{\text{архів}}}{\text{SIZE}_{\text{оригінал}}} \quad (1)$$

Таблиці результатів:

Табл. 1: Коефіцієнт стиснення для різних типів файлів

Тип файлів	Кількість файлів	$K_{\text{середнє}}$	$K_{\text{мін}}$	$K_{\text{макс}}$	Коментар
TXT	10	0.4570	0.4051	0.4969	добре стискаються
CSV	10	0.4324	0.1311	0.7861	добре стискаються
BMP	10	0.5120	0.0231	1.2182	добре стискаються
PDF	10	1.3300	1.2533	1.4080	розмір збільшився
EXE	10	0.6381	0.2554	0.9022	добре стискаються

PDF вже стиснутий у середині, а LZW — це ще один словниковий алгоритм, який не може вдруге стиснути вже стиснені дані. Тому результати стиснення коректні. Розмір збільшився через наявність хедера та кодів фіксовоної довжини (16 біт).

Табл. 2: Порівняння середнього коефіцієнту стиснення Хаффмана та LZW

Тип файлів	$K_{\text{середнє}} \text{Хаффман}$	$K_{\text{середнє}} \text{LZW}$
TXT	0.5917	0.4570
CSV	0.5965	0.4324
BMP	0.6817	0.5120
PDF	0.9955	1.3300
EXE	0.6839	0.6381

Висновки за результатами порівняння алгоритмів Хаффмана та LZW:

Порівняльний аналіз показав, що алгоритм LZW забезпечує кращу ефективність стиснення для текстових та табличних файлів порівняно з алгоритмом Хаффмана. Для бінарних і виконуваних файлів обидва алгоритми демонструють подібні результати. Для вже стиснених форматів, таких як PDF, застосування LZW є неефективним і призводить до збільшення розміру файлів.

4 Висновки

У результаті проведеного порівняльного аналізу ефективності стиснення даних алгоритмом LZW було встановлено, що ефективність стиснення суттєво залежить від типу вхідних даних та рівня їх повторюваності. Найкращі результати стиснення було отримано для табличних та текстових файлів. Для файлів формату CSV середній коефіцієнт стиснення становив 0.4324, а для текстових файлів TXT — 0.4570, що свідчить про високу ефективність алгоритму для даних з повторювальною структурою. Файли формату BMP та EXE продемонстрували середню ефективність стиснення. Для зображень BMP середнє значення коефіцієнта становило 0.5120, однак спостерігалася значна варіативність результатів, що пояснюється різною структурою зображень. Для виконуваних файлів EXE середній коефіцієнт стиснення склав 0.6381, що є типовим результатом для бінарних файлів зі складною внутрішньою структурою. Найгірші результати було зафіковано для файлів формату PDF, де середній коефіцієнт стиснення перевищив одиницю (1.3300). Це означає, що після застосування алгоритму LZW розмір файлів збільшився. Така поведінка пояснюється тим, що формат PDF зазвичай вже містить внутрішні механізми стиснення, а додаткове словниково кодування не дає позитивного ефекту. Отримані результати підтверджують, що алгоритм LZW найбільш ефективний для даних з високим рівнем повторюваності, зокрема текстових і табличних файлів, і малоекспективний або неекспективний для форматів, які вже містять стиснені дані. Загалом реалізація алгоритму LZW коректно працює для різних типів файлів та демонструє очікувані властивості словникового стиснення.