

Raport z laboratorium 3

Filip Nikolow

9 kwietnia 2021

1 Cel laboratorium

Celem laboratorium było zapoznanie się z algorytmami kompresji tekstu na przykładzie statycznego oraz dynamicznego algorytmu Huffmana. Jako zadanie dodatkowe należało zaimplementować dowolny algorytm ze zmiennym blokiem kompresji, który uzyska lepszy stopień kompresji od algorytmów Huffmana

2 Realizacja

Zaimplementowałem oba algorytmy Huffmana: statyczny oraz dynamiczny, a także w ramach zadania dodatkowego zaimplementowałem algorytm Lempela-Ziva-Welcha (LZW).

3 Realizacja poszczególnych poleceń

3.1 Polecenie 1 - format pliku

3.1.1 Statyczne kodowanie Huffmana

Zakodowany plik rozpoczyna się trzema bitami reprezentującymi liczbę zer (0-7) dopisanych na końcu pliku w ramach dopełnienia do pełnego bajtu. Następnie zakodowane jest drzewo Huffmana jako jego przeszukiwanie post-order:

- bit 0 oznacza węzeł wewnętrzny;
- bit 1 oznacza liść zawierający kodowany znak, kolejne 8 bitów to reprezentacja tego znaku w kodzie ASCII

Kodowanie drzewa jest zakończone pojedynczym bitem 0. Łatwo policzyć, że rozmiar reprezentacji liścia to 9 bitów, a węzła wewnętrznego 1 bit. Jak wiadomo, drzewo huffmana ma $2^m - 1$ węzłów (gdzie m to rozmiar alfabetu), a zatem uwzględniając kończący bit zerowy, tak zakodowane drzewo huffmana zajmuje dokładnie $10 \cdot m$ bitów.

Po drzewie następuje ciąg bitów reprezentujących zakodowany tekst.

3.1.2 Dynamiczne kodowanie Huffmana

Podobnie jak przy kodowaniu statycznym, pierwsze 3 bity reprezentują liczbę uzupełniających zer na końcu pliku. Po nich następuje ciąg bitów odpowiadających zakodowanemu dynamicznie tekstowi.

3.1.3 LZW

Ten algorytm nie potrzebuje żadnych dodatkowych metadanych w nagłówku pliku, gdyż pojedyncze bloki kodowania są nie mniejsze niż 8 bitów a więc dekodery jest w stanie samodzielnie odrzucić zera z końca pliku. Cały plik jest więc kodem zwracanym przez encoder.

3.2 Polecenie 2 - implementacja

Tak jak wspomniałem we wstępie, zaimplementowałem trzy algorytmy kompresji i dekompresji - 2 obowiązkowe i jeden dodatkowy. Z ważniejszych szczegółów implementacji algorytmu LZW: jako podstawowy zbiór znaków przyjąłem 256 znaków ASCII, a testy przeprowadziłem z limitem rozmiaru pojedynczego bloku kodującego ustawionym na 16 bitów (w efekcie kodowane bloki zmieniały się od 9 do 16) - można go łatwo zmienić, jest to stała klasy LZW. Początkowo testy przeprowadzałem na maksymalnym bloku równym 12 bitów, jednakże taki słownik był nieco za mały do efektywnego kodowania języka naturalnego - nadal algorytm LZW wygrywał większość testów, lecz ze znacznie mniejszą przewagą w przypadku kompresji książki.

3.3 Polecenia 3-4

Do testów pobrałem "Krzyżaków" w wersji angielskiej z portalu Guttenberga, oba zasugerowane na ćwiczeniach pliki źródłowe w języku C, oraz wygenerowałem za pomocą polecenia `base64` plik z losowymi znakami `ascii`. Pliki źródłowe połączyłem w jeden, a następnie je oraz książkę przepuściłem przez skrypt pythonowy podmieniający znaki `unicode` na najbardziej zbliżone znaki `ASCII` - wszystkie zaimplementowane algorytmy działają tylko na znakach `ASCII`.

4 Benchmarki

Poniżej załączam wszystkie wymagane pomiary. Zmierzone czasy kompresji i dekompresji na wykresach po lewej stronie są wyrażone w sekundach (a zatem im mniej tym lepiej), natomiast współczynniki kompresji na wykresach po stronie prawej są obliczone wzorem:

$$1 - \frac{\text{plik_skompresowany}}{\text{plik_nieskompresowany}}$$

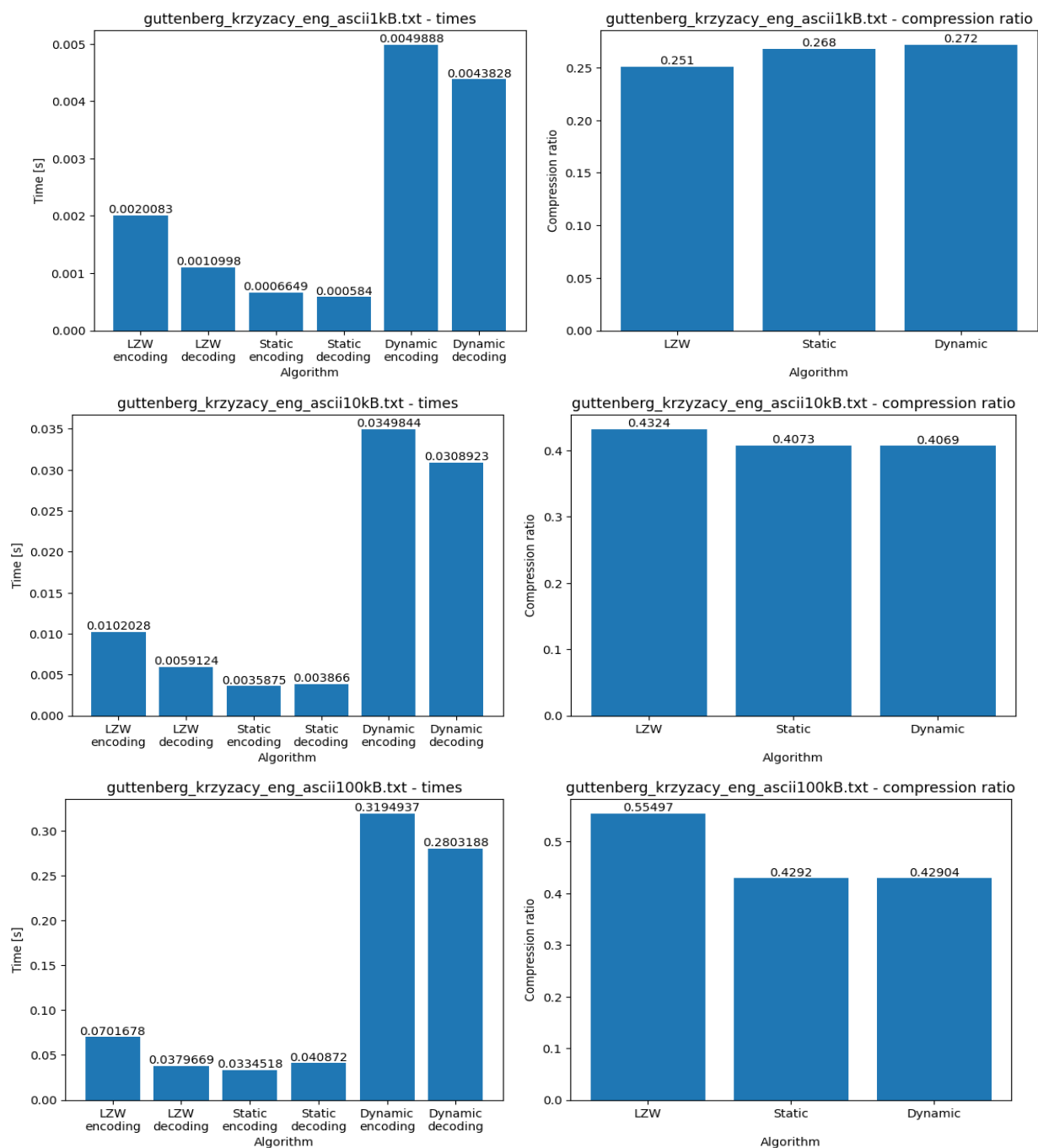
, a zatem im więcej tym lepiej. Porównując oba algorytmy Huffmana, zgodnie z oczekiwaniami algorytm dynamiczny osiąga minimalnie lepsze stopnie kompresji dla małych plików, natomiast przy plikach $\geq 10\text{KB}$ różnica się zaciera, a nawet zdarza się iż to kodowanie statyczne minimalnie wygrywa - są to jednak różnice raczej zaniedbywalne. Kodowanie dynamiczne za to zdecydowanie przegrywa jeśli chodzi o czasy kompresji i dekompresji - wina leży po stronie funkcji transformującej drzewo huffmana podczas czytania tekstu, gdyż operacje wyszukiwania węzłów do zamiany oraz same zamiany są kosztownymi operacjami. Przechodząc do algorytmu LZW także wielkich zaskoczeń nie ma: dla małych plików algorytmy Huffmana i LZW mają zbliżone osiągi, natomiast w miarę wzrostu rozmiaru pliku, ujawnia się też przewaga kodowania zmiennymi blokami nad kodowaniem blokami stałymi. W efekcie na testach z książką oraz kodem źródłowym, LZW przegrywa tylko raz: na najmniejszym fragmencie książki. Z kolei po kompresji tekstu z losowymi znakami nie ma co spodziewać się zbyt wiele - oczekiwanym rezultatem jest stopień kompresji bliski 0. Tak się jednak nie dzieje, gdyż algorytmy huffmana osiągają ok. 25% stopień kompresji, a nawet i algorytm LZW dla większych plików jest w stanie lekko zmniejszyć ich rozmiar. Jak się nad tym zastanowić, to także nie jest niczym zaskakującym, gdyż funkcja użyta do generowania znaków, `base64`, nie używa pełnej palety kodu `ASCII`, a więc kodowanie Huffmana skraca długość pliku nie poprzez podmienienie kodów najczęściej występujących liter na najkrótsze, a po prostu poprzez usunięcie z kodu niewystępujących znaków. Z kolei algorytm LZW ma minimalny blok kodowania równy 8 (a efektywnie 9), a zatem jest w stanie kompresować tylko powtarzające się wzorce - na małych plikach ilość takich powtórzeń jest na tyle niewielka, że ich kodowanie nie jest w stanie przewyższyć dużego narzutu na pojedyncze kody. Przy większych plikach słownik LZW zapełnia się i znowu, ponieważ `base64` generuje ograniczony alfabet, w końcu ilość powtarzających się wzorców jest na tyle duża, że możemy zaobserwować kompresję. Podsumowując, algorytm LZW wygrał 7 z 12 testów i w mojej opinii był najprostszym w kompletnej implementacji (wraz z formatem plików itp.). Bez

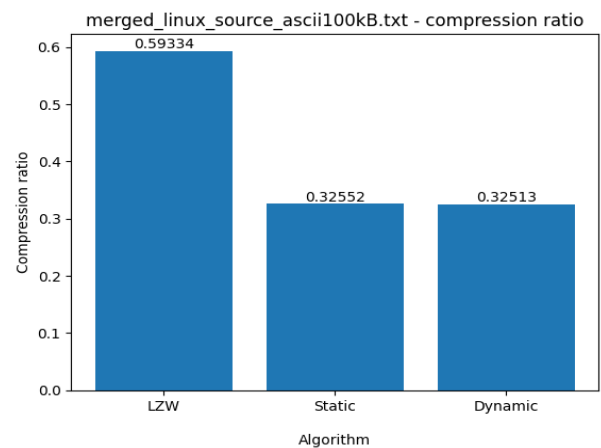
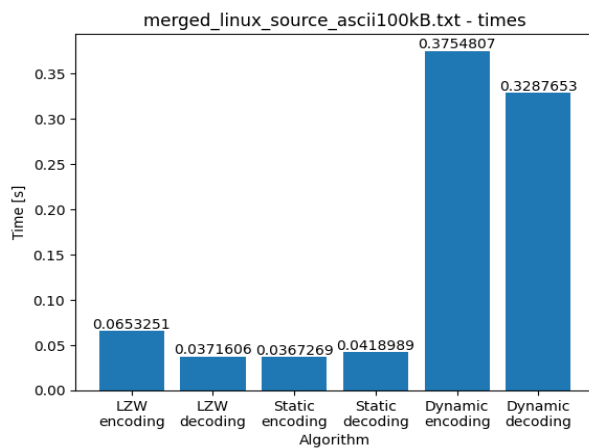
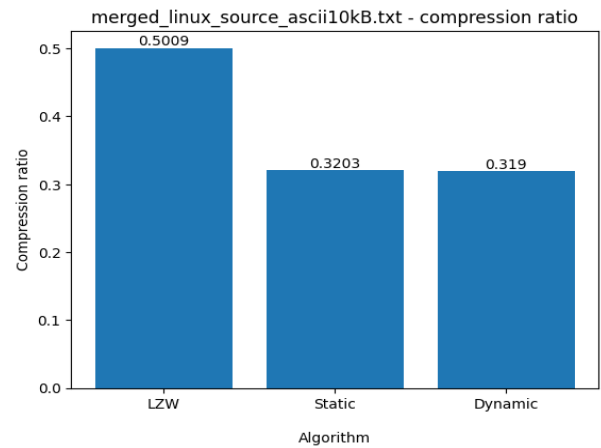
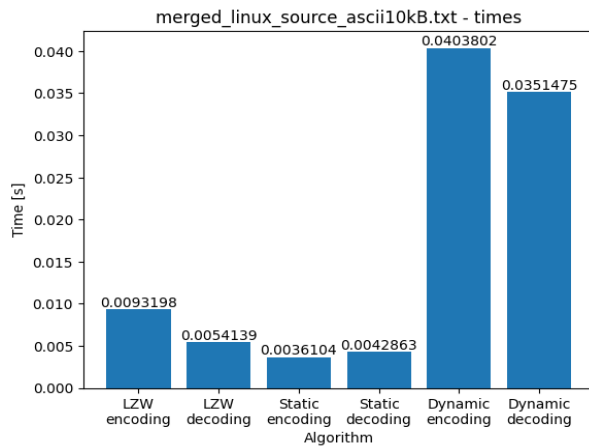
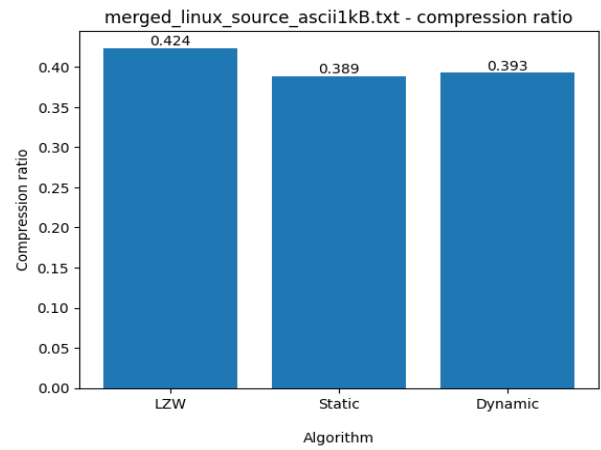
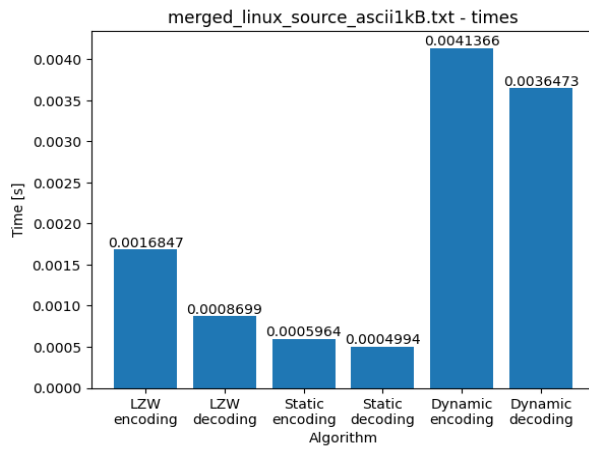
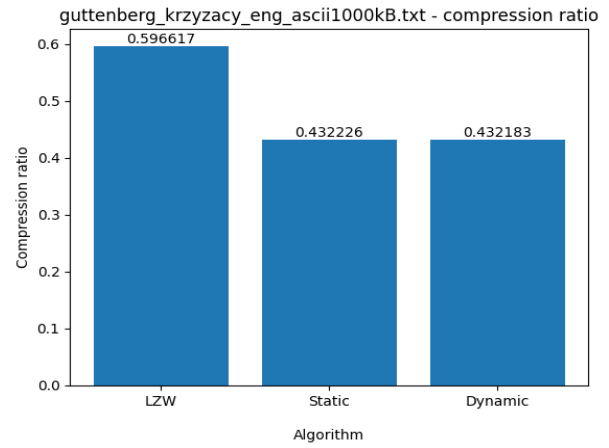
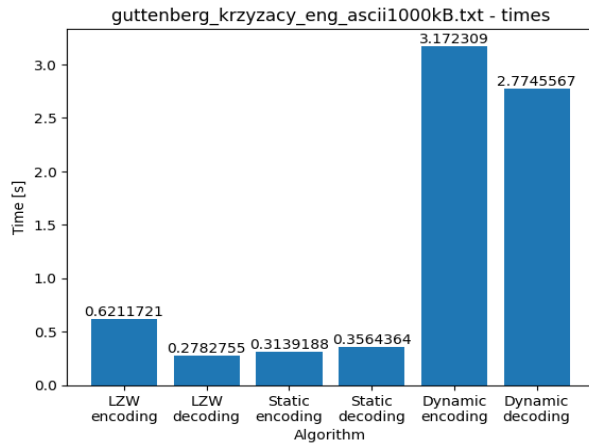
problemu mógłby wygrać 8 testów, gdyby jako podstawę kodowania przyjąć 128 znaków ASCII, a nie 256, ale byłby wtedy mniej uniwersalny.

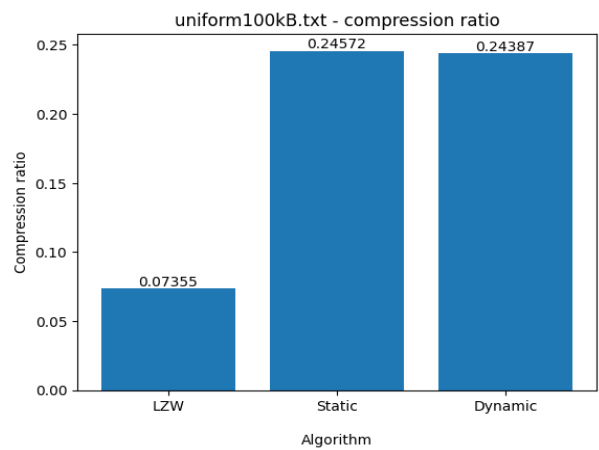
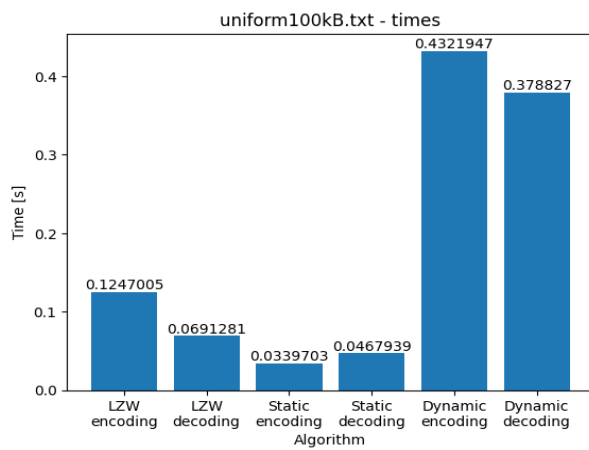
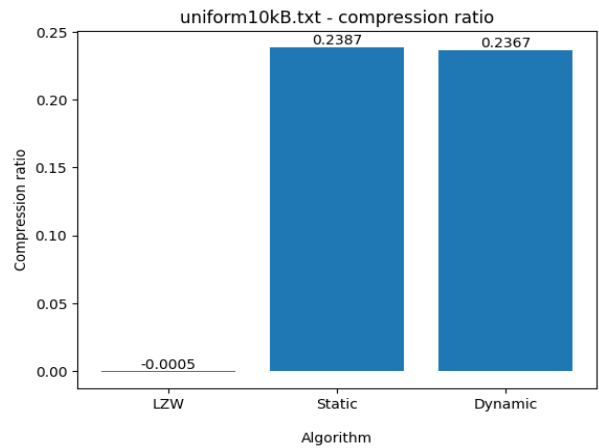
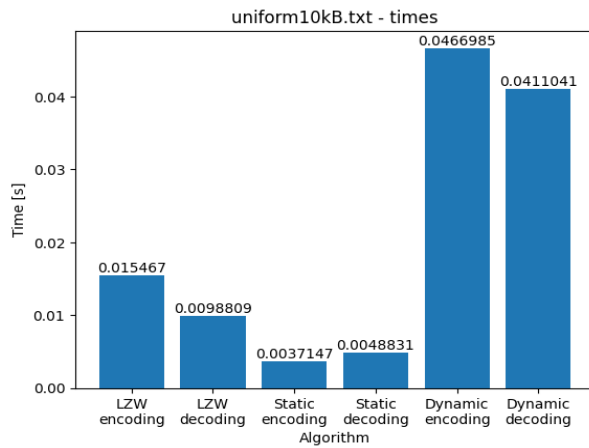
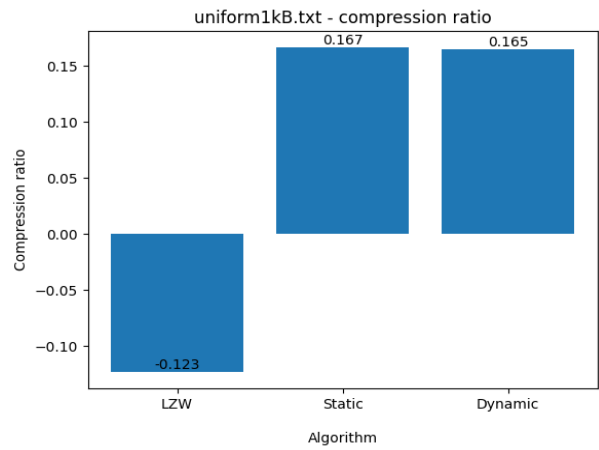
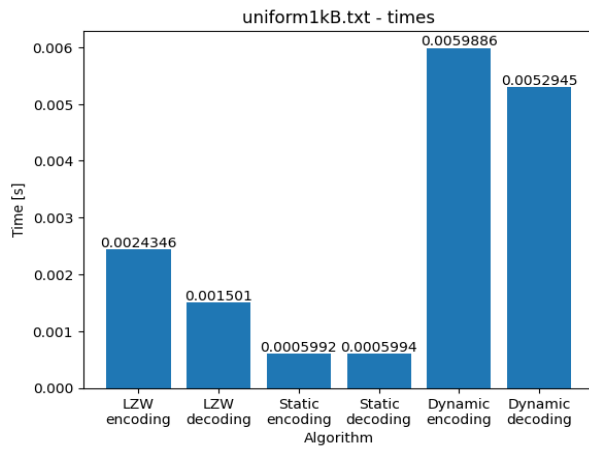
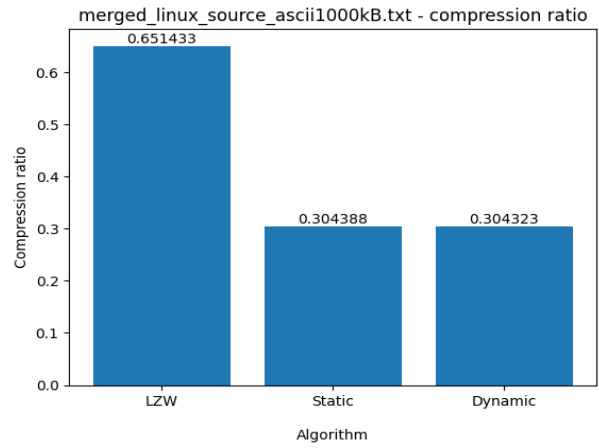
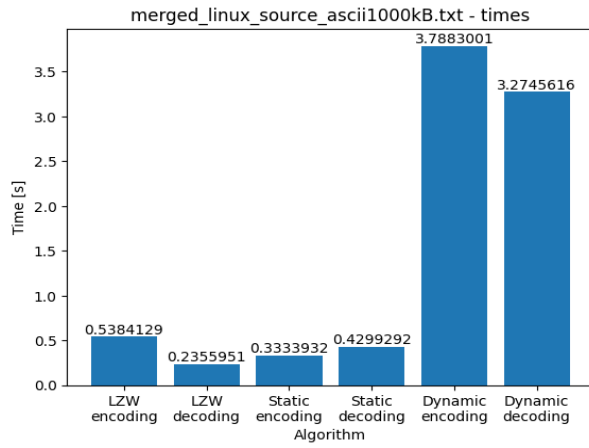
5 Szczegóły techniczne platformy testowej

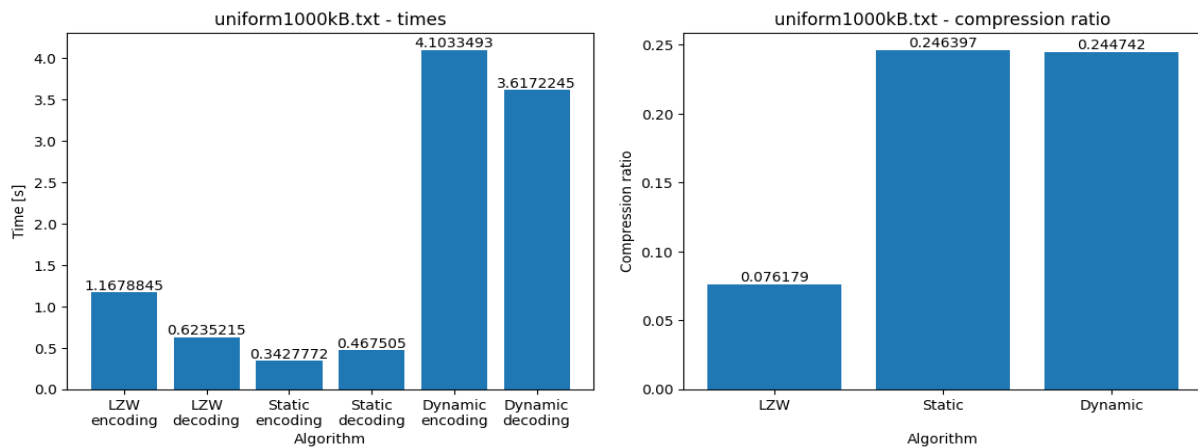
Wszystkie pomiary czasów zostały przeprowadzone na maszynie wirtualnej VirtualBox z systemem Ubuntu 20. Używana wersja języka Python to 3.8.5. Użyty procesor to AMD Ryzen 3900X. Poniższe pomiary czasów zostały powtórzone dziesięciokrotnie, a na wykresach zamieszczono średnie z tych pomiarów.

6 Benchmarki - wykresy









7 Kody modułów

7.1 Statyczne kodowanie Huffmana

```

1  from collections import deque
2  from bytearray import bytearray
3  from bytearray.util import ba2int, int2ba
4  from pretty_print import PrettyPrint
5  from node import Node
6
7
8  class StaticHuffman(PrettyPrint):
9
10     def __init__(self, text):
11         super().__init__()
12         self.text = text
13
14     def count_letter_frequency(self):
15         atom = 1 / len(self.text)
16         frequency = dict()
17         for c in self.text:
18             val = frequency.get(c, 0)
19             frequency[c] = val + atom
20         return frequency
21
22     def get_two_min(self, A, B):
23         res = []
24         while len(res) < 2:
25             a, b = 1, 1
26             if A:
27                 a = A[0].weight
28             if B:
29                 b = B[0].weight
30             if a < b:
31                 res.append(A.popleft())
32             else:
33                 res.append(B.popleft())
34         return res[0], res[1]

```

```

35
36 def static_huffman(self, letter_frequency):
37     trees = deque()
38     leaves = [Node(label, frequency) for label, frequency in letter_frequency.items()]
39     if len(leaves) == 1:
40         return leaves[0]
41     leaves.sort(key=lambda node: node.weight)
42     leaves = deque(leaves)
43     while leaves or len(trees) > 1:
44         N1, N2 = self.get_two_min(trees, leaves)
45         trees.append(N1.join(N2))
46     return trees[0]
47
48 def calculate_huffman_tree(self):
49     self.root = self.static_huffman(self.count_letter_frequency())
50     return self.root
51
52
53 class StaticCompressor:
54     ''' This class can take any class with a 'calculate_huffman_tree()' method
55     which returns a prefix-free tree root (using nodes from the 'Node' import)
56     Probably not a very useful feature :)'''
57
58     # File encoding metadata
59     TRAILING_ZEROS_SIZE = 3
60
61     def __init__(self, filename, huffman_class=None):
62         self.filename = filename
63         self.huffman_tree_builder = huffman_class
64
65     def encode(self, text, display_tree=False):
66         '''Encodes 'text' to a file named self.filename\n
67         current formatting is:
68         bits[:3] - trailing zeros count (automatically appended zeros at the end of file)
69         then the tree is encoded followed by a single 0 bit, then the text in huffman code
70         followed by 0-7 trailing zeros. Something like this:
71         [trailing_zeros, tree, '0', text]'''
72
73     def tree_encoder(node, code=''):
74         '''Does post-order traversal of the huffman tree
75         When it finds an inner node, '0' is printed
76         When it finds a leaf, it prints '1' followed by 8-bit ascii code
77         The tree size is 2m-1 where m is the alphabet size
78         hence the post-order traversal takes (2m-1) + 8*m = (10m - 1) bits of space'''
79         if node.left:
80             tree_encoder(node.left, code + '0')
81         if node.right:
82             tree_encoder(node.right, code + '1')
83         if node.label:
84             bits.append(1)
85             bits.frombytes(node.label.encode())
86             header_bit_count[0] += 9
87             ascii_to_bin[node.label] = code

```

```

88         else:
89             bits.append(0)
90             header_bit_count[0] += 1
91
92     bits = bytearray()
93     ascii_to_bin = dict()
94     header_bit_count = [self.TRAILING_ZEROS_SIZE]
95     H = self.huffman_tree_builder(text)
96     tree_encoder(H.calculate_huffman_tree())
97     bits.append(0)
98     if display_tree:
99         H.pretty_print()
100     for c in text:
101         bits.extend(bitarray(ascii_to_bin[c]))
102     tmp = bytearray()
103     length = (len(bits) + self.TRAILING_ZEROS_SIZE)
104     tmp.extend(int2ba((8 - (length % 8)) % 8,
105                     self.TRAILING_ZEROS_SIZE)) # trailing zeros count
106     tmp.extend(bits)
107     with open(self.filename, "wb") as f:
108         tmp.tofile(f)
109
110     def decode(self):
111         '''Decodes file 'self.filename' using formatting explained in the 'encode()' docstring'''
112
113     def tree_decoder():
114         stack = []
115         i = self.TRAILING_ZEROS_SIZE
116         while bits[i] is not False or len(stack) != 1:
117             if bits[i]:
118                 i += 1
119                 label = bits[i:i + 8].tobytes().decode()
120                 i += 8
121                 stack.append(Node(label))
122             else:
123                 i += 1
124                 right = stack.pop()
125                 left = stack.pop()
126                 stack.append(Node(left=left, right=right))
127         return stack[0], i + 1
128
129     bits = bytearray()
130     with open(self.filename, "rb") as f:
131         bits.fromfile(f)
132     trailing_zeros_count = ba2int(bits[:self.TRAILING_ZEROS_SIZE])
133     node, header_bit_count = tree_decoder()
134     root = node
135     # In case there is only one letter in the tree:
136     if node.left is None:
137         node.left = node
138         node.right = node
139     charlist = []
140     for i in range(header_bit_count, len(bits) - trailing_zeros_count):

```



```

141         if node.label is not None:
142             charlist.append(node.label)
143             node = root
144         if bits[i]:
145             node = node.right
146         else:
147             node = node.left
148
149     charlist.append(node.label)
150     return ''.join(charlist)
151
152
153 if __name__ == '__main__':
154     text = "abracadabra"
155     S = StaticCompressor('compressed.static', StaticHuffman)
156     S.encode(text, True)
157     print(S.decode())

```

7.2 Dynamiczne kodowanie Huffmana

```

1  from bitarray import bitarray
2  from bitarray.util import ba2int, int2ba
3  from pretty_print import PrettyPrint
4  from node import Node
5
6
7  class ExtendedNode(Node):
8
9      def __init__(self, index, label=None, weight=None, left=None, right=None, parent=None):
10         super().__init__(label, weight, left, right)
11         self.parent = parent
12         self.index = index
13
14     def swap_child_link(self, new_node, old_node):
15         if self.left == old_node:
16             self.left = new_node
17         else:
18             self.right = new_node
19
20     def swap_index(self, node, lis):
21         lis[self.index], lis[node.index] = lis[node.index], lis[self.index]
22         self.index, node.index = node.index, self.index
23
24     def swap(self, node, lis):
25         self.swap_index(node, lis)
26         if self.parent:
27             self.parent.swap_child_link(node, self)
28         if node.parent:
29             node.parent.swap_child_link(self, node)
30         node.parent, self.parent = self.parent, node.parent
31

```

```

32     def get_code(self):
33         code = []
34         node = self
35         while node.parent is not None:
36             if node.parent.left == node:
37                 code.append(False)
38             else:
39                 code.append(True)
40             node = node.parent
41         return reversed(code)
42
43
44     class DynamicHuffmanCompressor(PrettyPrint):
45         '''This class realizes dynamic huffman encoding and decoding'''
46
47         def __init__(self, filename):
48             super().__init__()
49             self.filename = filename
50             self.nodelist = []
51
52         def increment(self, node):
53             '''Takes care of weight increases and swaps if required'''
54             swap_index = node.index
55             while swap_index >= 0 and node.weight == self.nodelist[swap_index].weight:
56                 swap_index -= 1
57             to_swap = self.nodelist[swap_index + 1]
58             if to_swap != node and to_swap != node.parent:
59                 self.nodelist[swap_index + 1].swap(node, self.nodelist)
60             node.weight += 1
61             if node.parent:
62                 self.increment(node.parent)
63
64         def create_inner_node(self, empty_node, new_leaf):
65             '''Creates new inner node, used when a new letter is met'''
66             new_inner_node = ExtendedNode(len(self.nodelist),
67                                           None,
68                                           0,
69                                           left=empty_node,
70                                           right=new_leaf,
71                                           parent=empty_node.parent)
72             if empty_node.parent:
73                 empty_node.parent.swap_child_link(new_inner_node, empty_node)
74             new_leaf.parent = new_inner_node
75             empty_node.parent = new_inner_node
76             self.nodelist.append(new_inner_node)
77             new_inner_node.swap_index(empty_node, self.nodelist)
78             return new_inner_node
79
80         def add_letter(self, c, empty_node):
81             '''Adds a new letter to the tree'''
82             new_leaf = ExtendedNode(len(self.nodelist), c, 1)
83             self.nodelist.append(new_leaf)
84             new_inner_node = self.create_inner_node(empty_node, new_leaf)

```

```

85         self.increment(new_inner_node)
86     return new_leaf
87
88     def encode(self, text):
89         '''Dynamically encodes text as binary; every time a new letter is met, binary code of
90         the empty_node (representing all not yet met letters) is appended to binary code followed
91         by the ascii code of the letter; otherwise if a letter has already been in the tree, its
92         current code is appended and its weight is increased which also means that tree update
93         (increasing parents weights and swapping some nodes) may be required'''
94         empty_node = ExtendedNode(0, "##", 0)
95         nodes = {"##": empty_node}
96         self.nodelist = [empty_node]
97         bits = bitarray()
98         for c in text:
99             if c in nodes:
100                 bits.extend(nodes[c].get_code())
101                 self.increment(nodes[c])
102             else:
103                 bits.extend(empty_node.get_code())
104                 bits.frombytes(c.encode())
105                 nodes[c] = self.add_letter(c, empty_node)
106
107         self.root = self.nodelist[0]
108         trailing_zeros = (8 - ((len(bits) + 3) % 8)) % 8
109         tmp = bitarray(int2ba(trailing_zeros, 3))
110         tmp.extend(bits)
111         with open(self.filename, "wb") as f:
112             tmp.tofile(f)
113
114     def decode(self):
115         """Decodes binary file encoded by the 'encode(text)' method. Similarly to the encoding,
116         here the huffman tree is also dynamically adjusted while reading the binary file. Hence
117         it is not required to explicitly save huffman tree in the file thus allowing potentially
118         slightly better compression ratio"""
119         empty_node = ExtendedNode(0, "##", 0)
120         nodes = {"##": empty_node}
121         self.nodelist = [empty_node]
122         bits = bitarray()
123         with open(self.filename, "rb") as f:
124             bits.fromfile(f)
125         trailing_zeros = ba2int(bits[:3])
126         charlist = []
127         i = 3
128         while i < len(bits) - trailing_zeros:
129             node = self.nodelist[0]
130             # nodes always have either both children or none, so checking for one is sufficient
131             while node.left:
132                 if bits[i]:
133                     node = node.right
134                 else:
135                     node = node.left
136                 i += 1
137             if node == empty_node:

```

```

138         letter = bits[i:i + 8].tobytes().decode()
139         i += 8
140         nodes[letter] = self.add_letter(letter, empty_node)
141     else:
142         letter = node.label
143         self.increment(nodes[letter])
144         charlist.append(letter)
145     return ''.join(charlist)
146
147
148 if __name__ == '__main__':
149     text = "abracadabra"
150     D = DynamicHuffmanCompressor('compressed.dynamic')
151     D.encode(text)
152     print(D.decode())

```

7.3 Lempel-Ziv-Welch

```

1  from bitarray import bitarray
2  from bitarray.util import ba2int, int2ba
3  from math import log2, ceil
4
5
6  class LZW:
7      """Lempel-Ziv-Welch algorithm implementation with variable compression block
8      The compression block dynamically changes as the dictionary fills. Default variability
9      range is from 9 to 16 (first character is encoded with 8 bits).
10     Changing block size can yield much better compression ratio,
11     especially on natural language texts"""
12
13     BLOCK_SIZE = 16
14     ASCII_LEN = 256
15
16     def __init__(self, filename):
17         self.filename = filename
18         self.MAX_LEN = int(2**self.BLOCK_SIZE)
19
20     def encode(self, text):
21         base_table = dict()
22         for i in range(self.ASCII_LEN):
23             base_table[chr(i)] = int2ba(i, self.BLOCK_SIZE)
24         bits = bitarray()
25         table = base_table.copy()
26         table_len = self.ASCII_LEN
27         string = text[0]
28         for i in range(1, len(text)):
29             char = text[i]
30             if string + char in table:
31                 string += char
32             else:
33                 bits.extend(table[string][-ceil(log2(table_len)):])

```

```

34         if table_len == self.MAX_LEN:
35             table = base_table.copy()
36             table_len = self.ASCII_LEN
37             table[string + char] = int2ba(table_len, self.BLOCK_SIZE)
38             table_len += 1
39             string = char
40         bits.extend(table[string][-ceil(log2(table_len)):])
41         with open(self.filename, "wb") as f:
42             bits.tofile(f)
43
44     def decode(self):
45         table = [None] * self.MAX_LEN
46         for i in range(self.ASCII_LEN):
47             table[i] = chr(i)
48         bits = bitarray()
49         with open(self.filename, "rb") as f:
50             bits.fromfile(f)
51         table_len = self.ASCII_LEN
52         start = ceil(log2(table_len))
53         string = table[ba2int(bits[:start])]
54         text = string
55         # If the character count is not a multiply of BLOCK_SIZE, there were some
56         # automatically appended zeros at the end of file to fill whole byte, we skip these
57         # It's not actually needed to work, but is a nice thing to expose
58         actual_size = len(bits) - (len(bits) % self.BLOCK_SIZE)
59         i = start
60         while i < actual_size:
61             if table_len == self.MAX_LEN:
62                 table_len = self.ASCII_LEN
63                 jump = ceil(log2(table_len + 1))
64                 code = ba2int(bits[i:i + jump])
65                 if code >= table_len:
66                     entry = string + string[0]
67                 else:
68                     entry = table[code]
69                 text += entry
70                 table[table_len] = string + entry[0]
71                 table_len += 1
72                 string = entry
73                 i += jump
74         return text
75
76
77 if __name__ == '__main__':
78     text = 'abracadabra'
79     C = LZW('compressed.lzw')
80     # C.encode(text)
81     print(C.decode())

```

7.4 Moduły pomocnicze

7.4.1 Node

```
1 class Node:
2     node_count = 0
3
4     def __init__(self, label=None, weight=None, left=None, right=None):
5         self.ind = Node.node_count
6         Node.node_count += 1
7         self.label = label
8         self.weight = weight
9         self.left = left
10        self.right = right
11
12    def join(self, node):
13        return Node(None, self.weight + node.weight, self, node)
14
15    def __repr__(self):
16        if self.weight:
17            return "Index: " + str(self.ind) + "\nLabel: " + str(
18                self.label) + "\nWeight: " + str(round(self.weight, 2))
19        else:
20            return "Index: " + str(self.ind) + "\nLabel: " + str(self.label)
```

7.4.2 Pretty Print

```
1 import pydot
2 import tempfile
3 from PIL import Image
4
5
6 class PrettyPrint:
7     '''Base class for printing trees using pydot's write method
8     Tree node needs to have 'left' and 'right' attributes'''
9
10    def __init__(self):
11        self.root = None
12
13    def pretty_print(self, name=None, display=True):
14        '''Prints a tree as an image\n
15        name - if provided, saves tree image to that filename,
16        display - wheter to display the image'''
17
18        def dfs_helper(node, graph):
19            for c in [node.left, node.right]:
20                if c:
21                    graph.add_edge(pydot.Edge(str(node), str(c)))
22                    dfs_helper(c, graph)
23
24        if self.root is None:
25            print("Root is None!")
```

```

26         return
27
28     graph = pydot.Dot(graph_type='graph')
29     dfs_helper(self.root, graph)
30     if name is None:
31         fout = tempfile.NamedTemporaryFile(suffix=".png")
32         name = fout.name
33     else:
34         name += "_tree.png"
35     graph.write(name, format="png")
36     if display:
37         Image.open(name).show()

```

7.4.3 Benchmark

```

1  from lempel_ziv_welch import LZW
2  from static_huffman import StaticHuffman, StaticCompressor
3  from dynamic_huffman import DynamicHuffmanCompressor
4  from time import time
5  import matplotlib.pyplot as plt
6  from os.path import getsize
7
8
9  def benchmark(funclist, datlist, labels, repeats=None):
10     if repeats is None:
11         repeats = 1
12     xs = []
13     ys = []
14     for dat in datlist:
15         times = dict()
16         for func, label in zip(funclist, labels):
17             for _ in range(repeats):
18                 t_enc, t_dec = func(dat)
19                 key1, key2 = label + '\nencoding', label + '\ndecoding'
20                 times[key1] = times.get(key1, 0) + (t_enc / repeats)
21                 times[key2] = times.get(key2, 0) + (t_dec / repeats)
22     xs.append(times.keys())
23     ys.append(times.values())
24     return xs, ys
25
26
27 def compression_ratio_benchmark(funclist, datlist, labels, repeats=1):
28     xs = []
29     ys = []
30     for dat in datlist:
31         times = dict()
32         for func, label in zip(funclist, labels):
33             for _ in range(repeats):
34                 ratio = func(dat)
35                 key = label + '\n'
36                 times[key] = times.get(key, 0) + (ratio) / repeats

```

```

37         xs.append(times.keys())
38         ys.append(times.values())
39     return xs, ys
40
41
42 def plot_factory(xs, ys, titles, axis_labels=None, save_to_file=False, display=True):
43     for i in range(len(xs)):
44         plt.bar(xs[i], ys[i])
45         plt.title(titles[i])
46         for j, v in enumerate(ys[i]):
47             plt.annotate(str(round(v, 7)), xy=(j, v), ha='center', va='bottom')
48         if axis_labels is not None:
49             plt.xlabel(axis_labels['x'])
50             plt.ylabel(axis_labels['y'])
51
52         if save_to_file:
53             plt.savefig("plots/" + titles[i] + ".png", bbox_inches='tight')
54         if display:
55             plt.show()
56
57
58 def universal_compression_wrapper(source_filename, compressed_filename, C):
59     uncompressed = getsize('../sources/' + source_filename)
60     with open('../sources/' + source_filename, "r") as f:
61         text = f.read()
62     C.encode(text)
63     compressed = getsize(compressed_filename)
64     return 1 - (compressed / uncompressed)
65
66
67 def LZW_compression_benchmark_wrapper(source_filename):
68     filename = 'compressed.lzw'
69     C = LZW(filename)
70     return universal_compression_wrapper(source_filename, filename, C)
71
72
73 def static_compression_benchmark_wrapper(source_filename):
74     filename = 'compressed.static'
75     C = StaticCompressor(filename, StaticHuffman)
76     return universal_compression_wrapper(source_filename, filename, C)
77
78
79 def dynamic_compression_benchmark_wrapper(source_filename):
80     filename = 'compressed.dynamic'
81     C = DynamicHuffmanCompressor(filename)
82     return universal_compression_wrapper(source_filename, filename, C)
83
84
85 def universal_wrapper(source_filename, C):
86     with open('../sources/' + source_filename, "r") as f:
87         text = f.read()
88     t_enc1 = time()
89     C.encode(text)

```



```

90     t_enc2 = time()
91     t_dec1 = time()
92     C.decode()
93     t_dec2 = time()
94     return t_enc2 - t_enc1, t_dec2 - t_dec1
95
96
97 def LZW_benchmark_wrapper(source_filename):
98     C = LZW('compressed.lzw')
99     return universal_wrapper(source_filename, C)
100
101
102 def static_benchmark_wrapper(source_filename):
103     C = StaticCompressor('compressed.static', StaticHuffman)
104     return universal_wrapper(source_filename, C)
105
106
107 def dynamic_benchmark_wrapper(source_filename):
108     C = DynamicHuffmanCompressor('compressed.dynamic')
109     return universal_wrapper(source_filename, C)
110
111
112 if __name__ == '__main__':
113     labels = ["LZW", "Static", "Dynamic"]
114     datlist = [
115         name + str(i) + 'kB.txt'
116         for i in [1, 10, 100, 1000]
117         for name in ['gutenberg_krzyzacy_eng_ascii', 'merged_linux_source_ascii', 'uniform']
118     ]
119     # time benchmarks
120     funclist = [LZW_benchmark_wrapper, static_benchmark_wrapper, dynamic_benchmark_wrapper]
121     xs, ys = benchmark(funclist, datlist, labels, 10)
122     # compression ratio benchmarks
123     funclist = [
124         LZW_compression_benchmark_wrapper, static_compression_benchmark_wrapper,
125         dynamic_compression_benchmark_wrapper
126     ]
127     xs2, ys2 = compression_ratio_benchmark(funclist, datlist, labels)
128     # Plotting
129     titles = [s + " - times" for s in datlist]
130     plot_factory(xs, ys, titles, {'x': 'Algorithm', 'y': 'Time [s]'}, True)
131     titles = [s + " - compression ratio" for s in datlist]
132     plot_factory(xs2, ys2, titles, {'x': 'Algorithm', 'y': 'Compression ratio'}, True)

```
