

НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ
«ВЫСШАЯ ШКОЛА ЭКОНОМИКИ»

ОТЧЕТ ПО ИНДИВИДУАЛЬНОМУ ДОМАШНЕМУ
ЗАДАНИЮ №2

**«Практический анализ алгоритмов
поиска вхождений строки-шаблона в тексте.»**

Выполнил:
студент 213 группы
Мовшин М. А.

Москва
2023

Содержание

1. Цель работы	2
2. Тестовые данные	3
3. Особенности проведения исследования	4
4. Исследование	5
4.1 Наивный алгоритм	5
4.2 КМП алгоритмы	12
4.3 Замечания	19
5. Выводы	20
6. Листинги	21

1. Цель работы

- Реализация различных алгоритмов поиска вхождений заданного паттерна (шаблона) в исходном тексте
- Экспериментальный анализ временной сложности алгоритмов
 - алгоритм Кнута-Морриса-Прата
 - алгоритм Кнута-Морриса-Прата с применением уточненных границ
 - наивный алгоритм
- Сопоставление теоретических и практических оценок временной сложности алгоритмов
- Интерпретация экспериментальных результатов

2. Тестовые данные

В соответствии с запросами условия мною были проведены тесты всех трех алгоритмов со всеми возможными вариациями следующих параметров:

- Длина текста: 10.000/100.000 символов
- Размер алфавита: 2/4 символов
- Количество символов подстановки: 1/2/3/4 символа

И размером паттерна от 100 до 3000 символов с шагом в 100 символов. Эти тесты, однако, не позволяли построить репрезентативные графики для КМП-алгоритмов - те отрабатывали слишком быстро, либо прирост размера паттерна был непоказателен. Поэтому я также провел тесты с теми же вариациями по размеру алфавита и количеству символов подстановки, но другими размерами текстов и паттернов.

Размер текста: 1.000.000 символов. Размер паттерна - до половины размера текста с шагом в $\frac{1}{50}$ размера текста.

3. Особенности проведения исследования

Для усреднения замеров за время прохождения теста бралось среднее по пяти запускам.

Для проверки паттерна с учетом символа подстановки все сравнения символов производились с помощью специальной функции, верно трактующей символ подстановки.

Генерация тестовых данных и выполнение тестов реализованы внутри одной программы. Время прохождения тестов не учитывает время генерации данных.

Результирующие данные были экспортированы в текстовые файлы. После они были загружены в среду Jupiter, где была проделана аналитическая работа - построение графиков и т. д.

4. Исследование

Для усреднения замеров за время прохождения теста бралось среднее по пяти запускам.

Для проверки паттерна с учетом символа подстановки все сравнения символов производились с помощью специальной функции, верно трактующей символ подстановки.

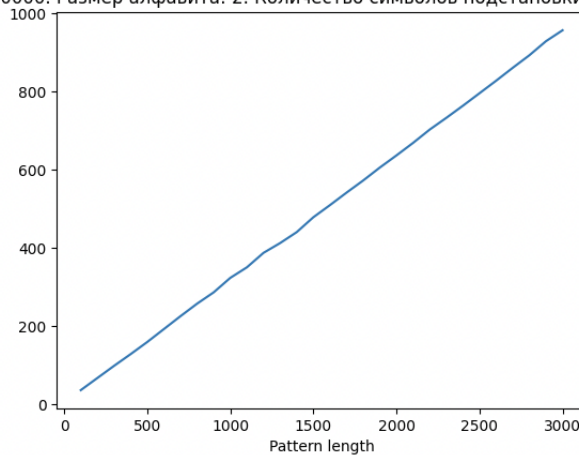
4.1 Наивный алгоритм

Ожидаемая асимптотика для наивного алгоритма:

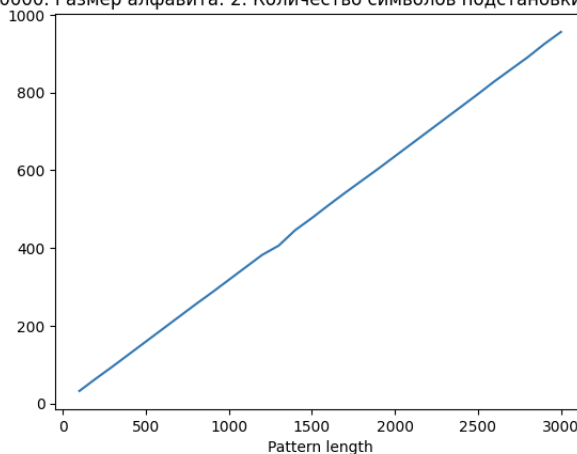
$$O(text_size \cdot pattern_size)$$

Взглянем на графики:

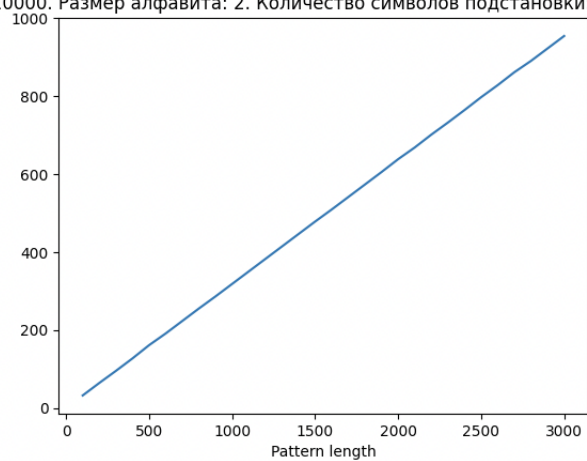
Размер текста: 10000. Размер алфавита: 2. Количество символов подстановки: 0. Наивный алгоритм.



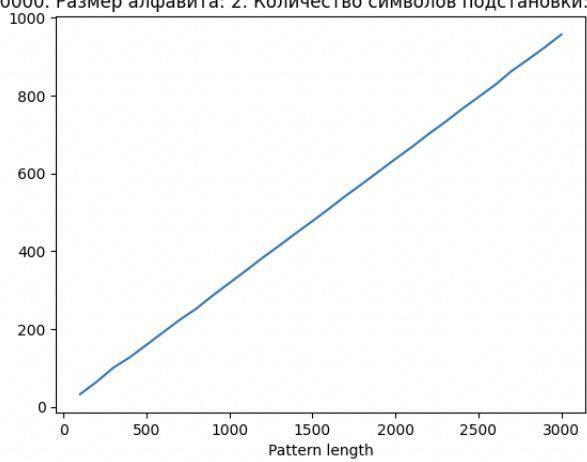
Размер текста: 10000. Размер алфавита: 2. Количество символов подстановки: 1. Наивный алгоритм.



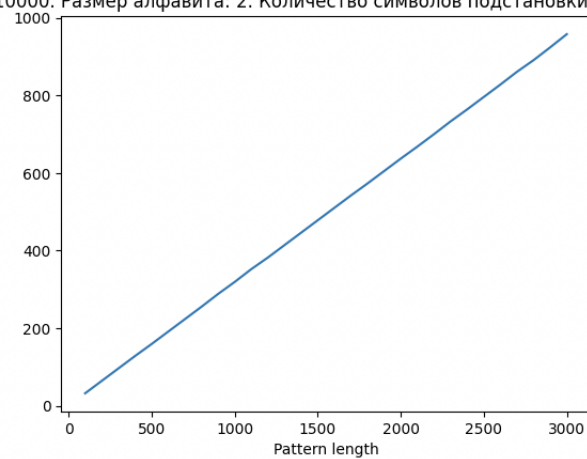
Размер текста: 10000. Размер алфавита: 2. Количество символов подстановки: 2. Наивный алгоритм.



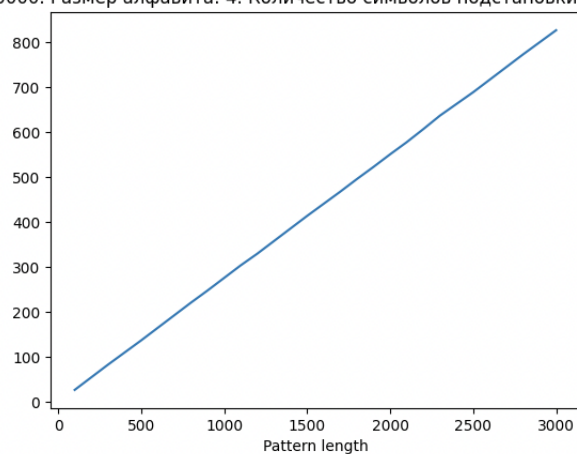
Размер текста: 10000. Размер алфавита: 2. Количество символов подстановки: 3. Наивный алгоритм.



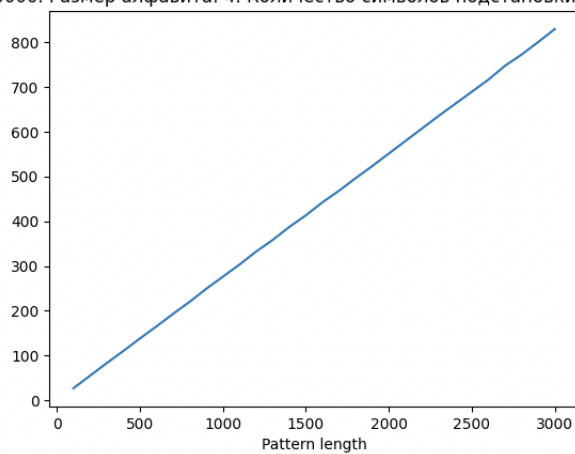
Размер текста: 10000. Размер алфавита: 2. Количество символов подстановки: 4. Наивный алгоритм.



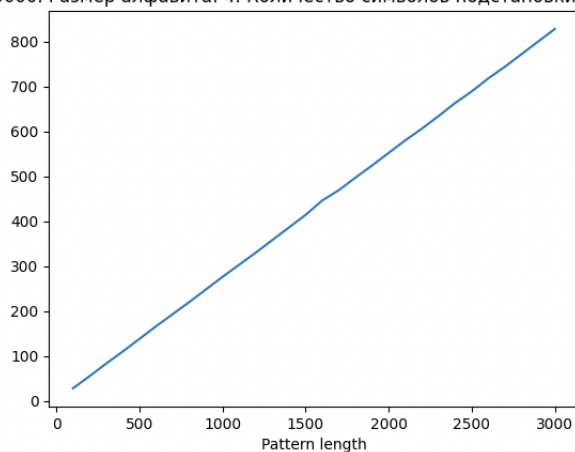
Размер текста: 10000. Размер алфавита: 4. Количество символов подстановки: 0. Наивный алгоритм.



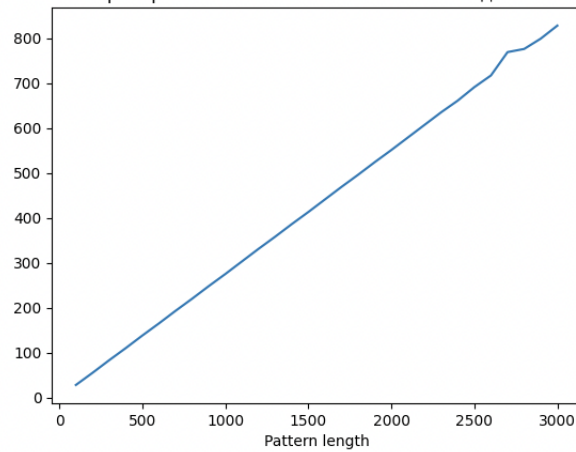
Размер текста: 10000. Размер алфавита: 4. Количество символов подстановки: 1. Наивный алгоритм.



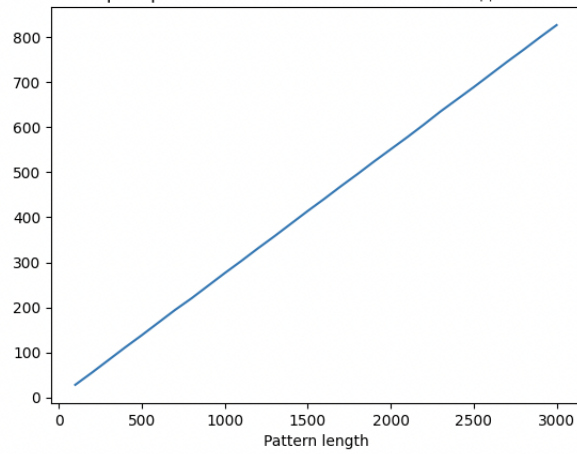
Размер текста: 10000. Размер алфавита: 4. Количество символов подстановки: 2. Наивный алгоритм.



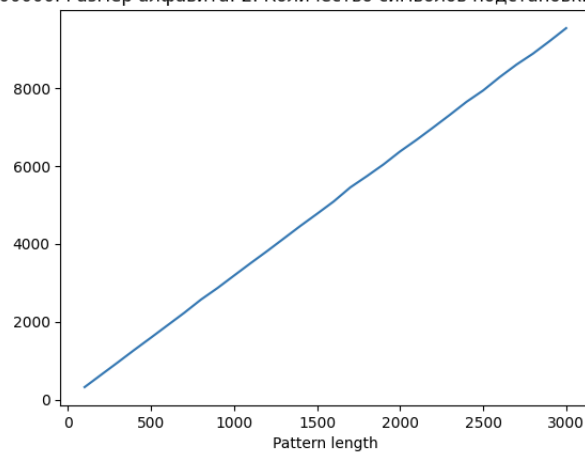
Размер текста: 10000. Размер алфавита: 4. Количество символов подстановки: 3. Наивный алгоритм.



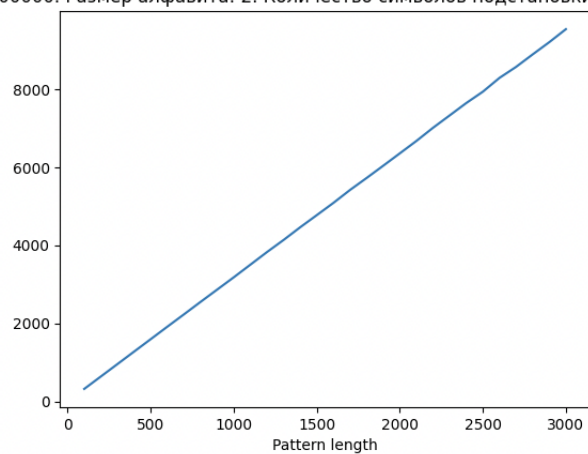
Размер текста: 10000. Размер алфавита: 4. Количество символов подстановки: 4. Наивный алгоритм.



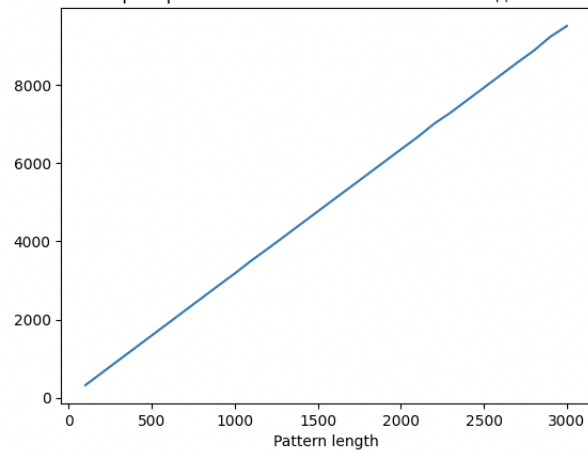
Размер текста: 100000. Размер алфавита: 2. Количество символов подстановки: 0. Наивный алгоритм.



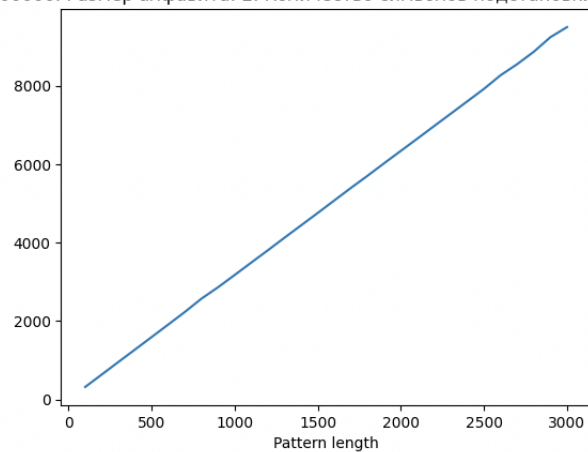
Размер текста: 100000. Размер алфавита: 2. Количество символов подстановки: 1. Наивный алгоритм.



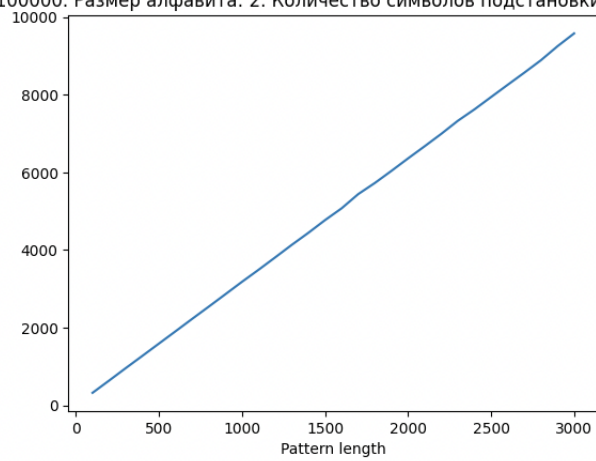
Размер текста: 100000. Размер алфавита: 2. Количество символов подстановки: 2. Наивный алгоритм.



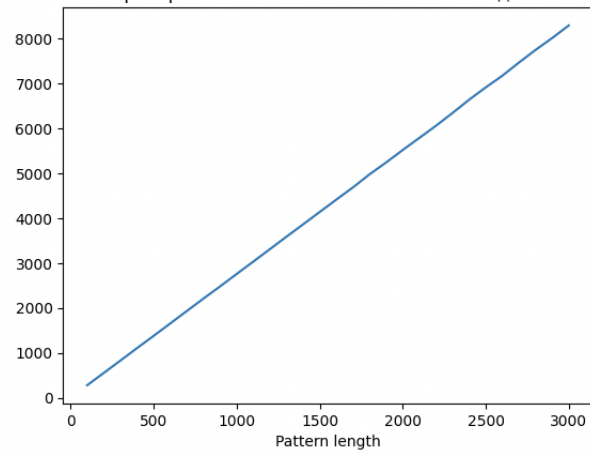
Размер текста: 100000. Размер алфавита: 2. Количество символов подстановки: 3. Наивный алгоритм.



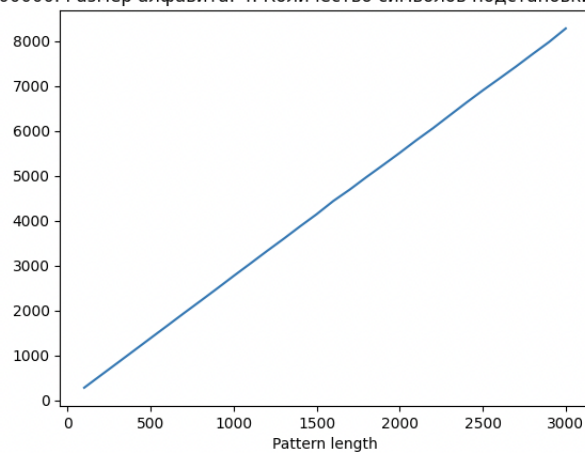
Размер текста: 100000. Размер алфавита: 2. Количество символов подстановки: 4. Наивный алгоритм.



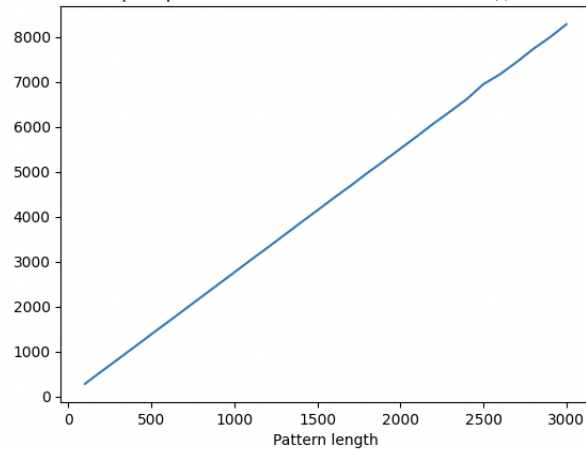
Размер текста: 100000. Размер алфавита: 4. Количество символов подстановки: 0. Наивный алгоритм.



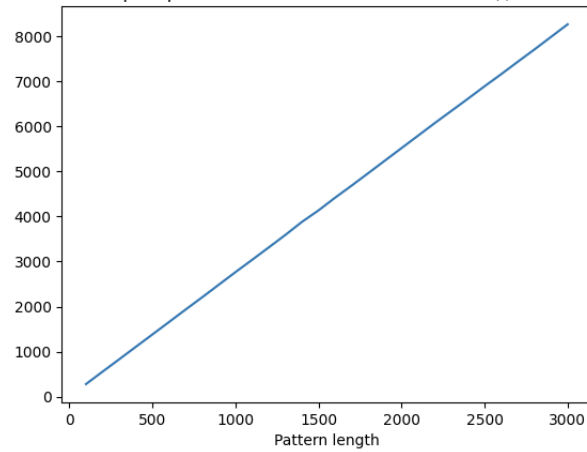
Размер текста: 100000. Размер алфавита: 4. Количество символов подстановки: 1. Наивный алгоритм.



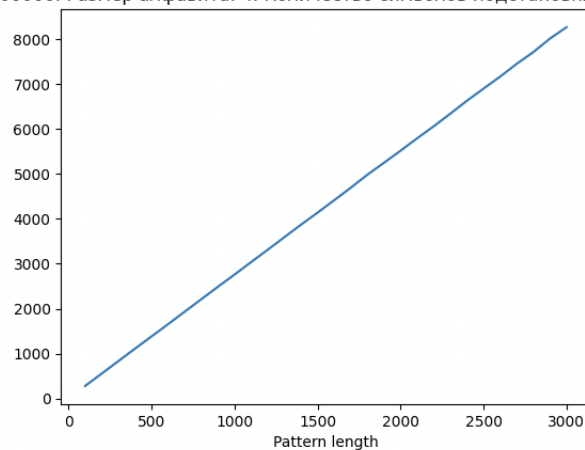
Размер текста: 100000. Размер алфавита: 4. Количество символов подстановки: 2. Наивный алгоритм.



Размер текста: 100000. Размер алфавита: 4. Количество символов подстановки: 3. Наивный алгоритм.



Размер текста: 100000. Размер алфавита: 4. Количество символов подстановки: 4. Наивный алгоритм.



На всех графиках прослеживается одна и та же картина - почти идеальная прямая стремящаяся началом к началу координат. Во-первых, легко заметить линейную зависимость от размера паттерна. Во-вторых, то, что прямая исходит из начала координат говорит о том, что размер паттерна выступает множителем в асимптотической формуле. По предложенным в условии входным данным

я сделал вывод, что линейная зависимость от размера текста считается очевидной (не было предложено делать графики с константным размером паттерна и изменяющимся размером текста). Поэтому можно сделать общий вывод об асимптотике: раз размер паттерна входит как множитель в первой степени, асимптотика для наивного алгоритма действительно:

$$O(text_size \cdot pattern_size)$$

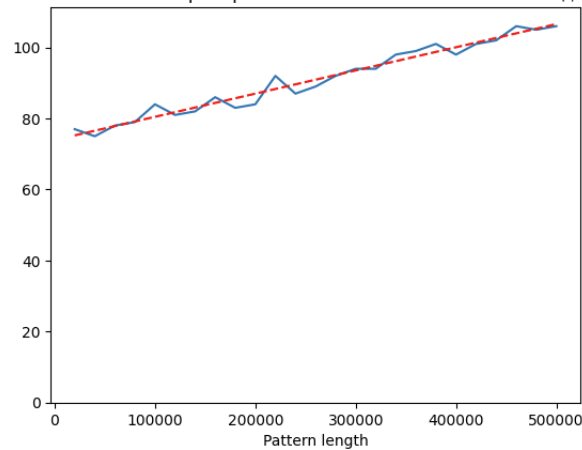
4.2 КМП алгоритмы

Ожидаемая асимптотика для КМП алгоритмов (обоих):

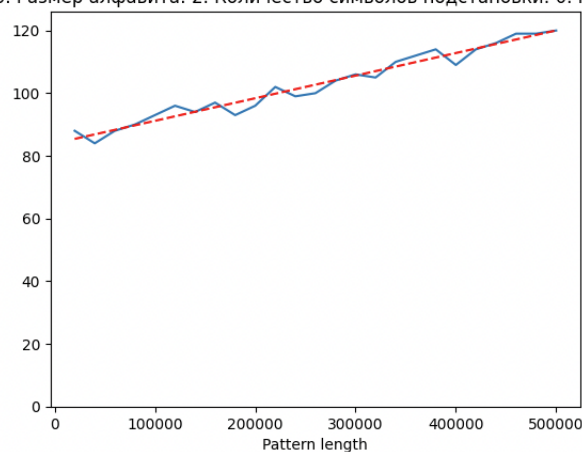
$$O(text_size + pattern_size)$$

Взглянем на графики:

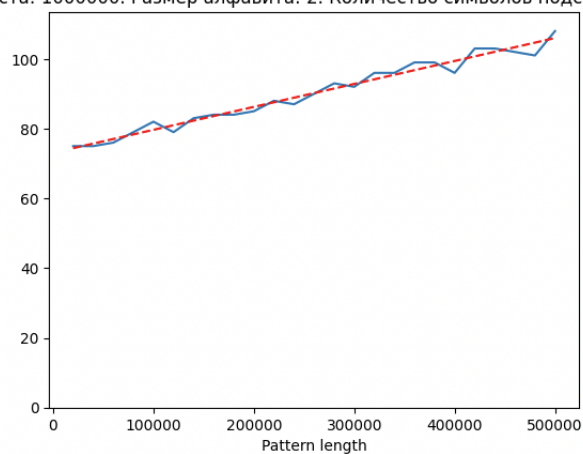
Размер текста: 1000000. Размер алфавита: 2. Количество символов подстановки: 0. КМП



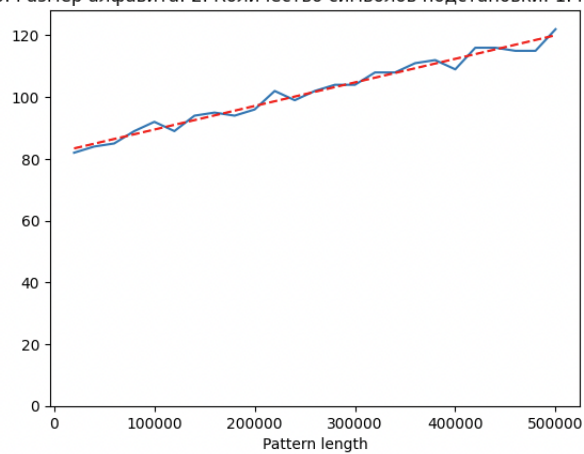
Размер текста: 1000000. Размер алфавита: 2. Количество символов подстановки: 0. КМП с уточненными границами



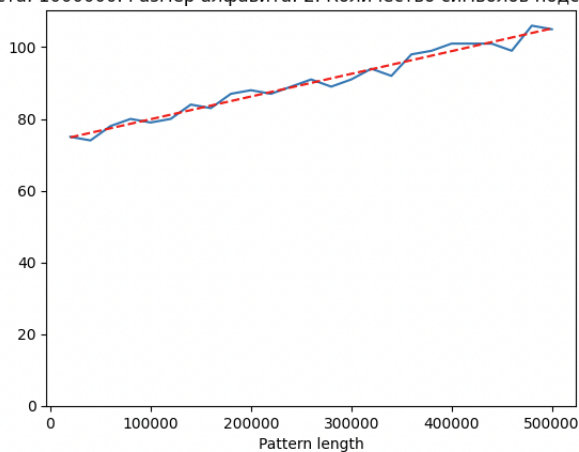
Размер текста: 1000000. Размер алфавита: 2. Количество символов подстановки: 1. КМП



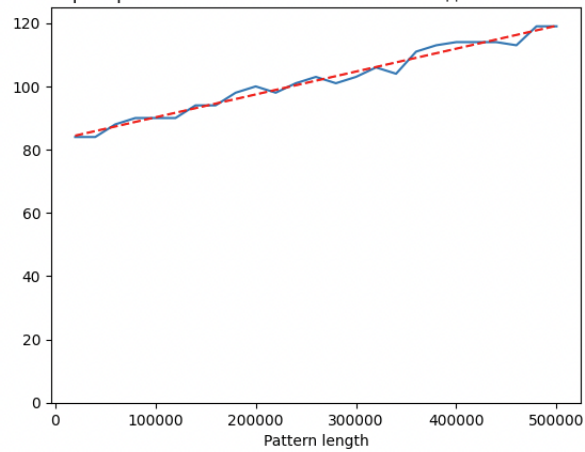
Размер текста: 1000000. Размер алфавита: 2. Количество символов подстановки: 1. КМП с уточненными границами



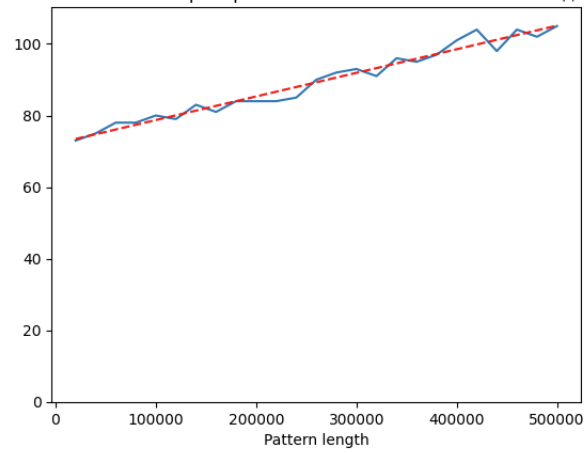
Размер текста: 1000000. Размер алфавита: 2. Количество символов подстановки: 2. КМП



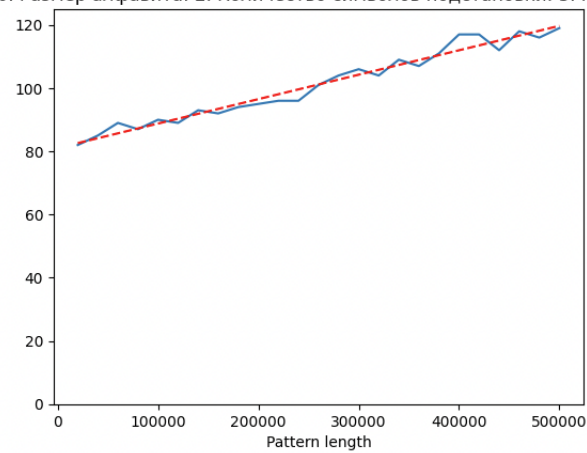
Размер текста: 1000000. Размер алфавита: 2. Количество символов подстановки: 2. КМП с уточненными границами



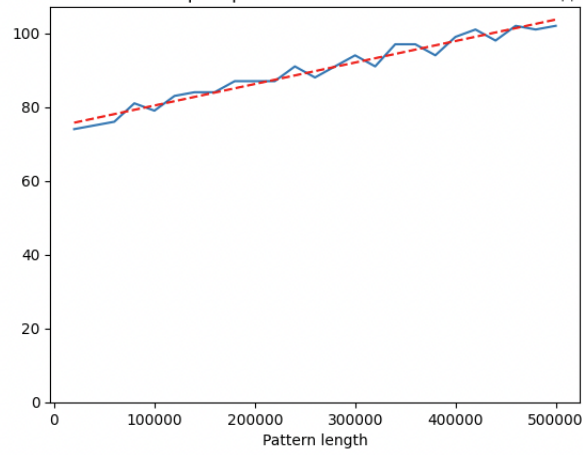
Размер текста: 1000000. Размер алфавита: 2. Количество символов подстановки: 3. КМП



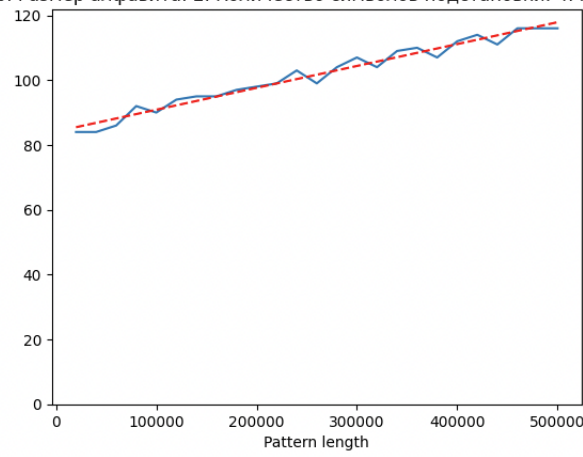
Размер текста: 1000000. Размер алфавита: 2. Количество символов подстановки: 3. КМП с уточненными границами



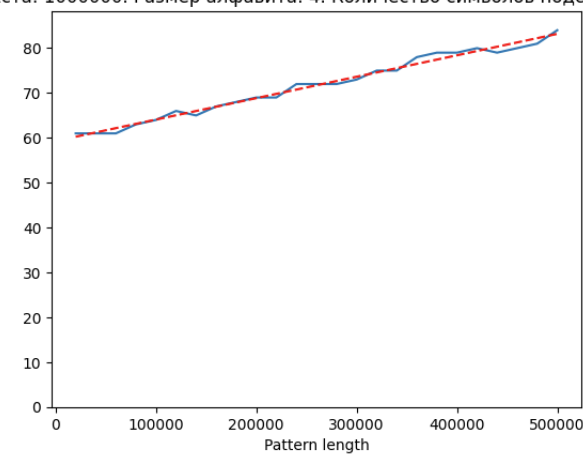
Размер текста: 1000000. Размер алфавита: 2. Количество символов подстановки: 4. КМП



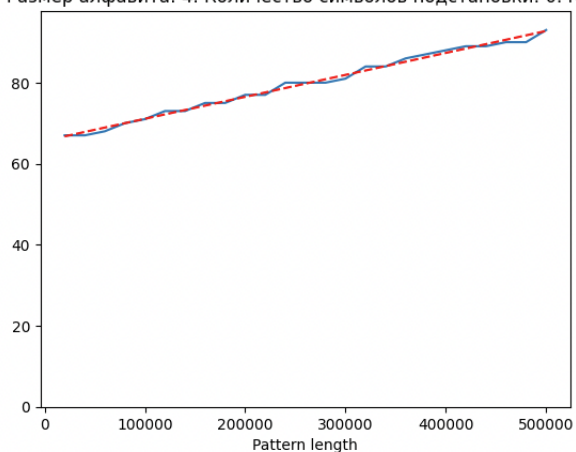
Размер текста: 1000000. Размер алфавита: 2. Количество символов подстановки: 4. КМП с уточненными границами



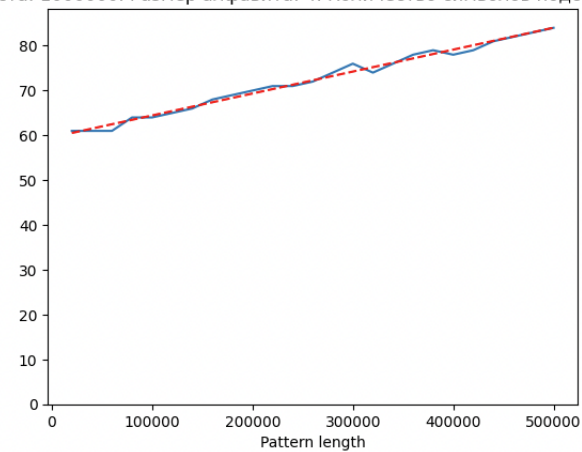
Размер текста: 1000000. Размер алфавита: 4. Количество символов подстановки: 0. КМП



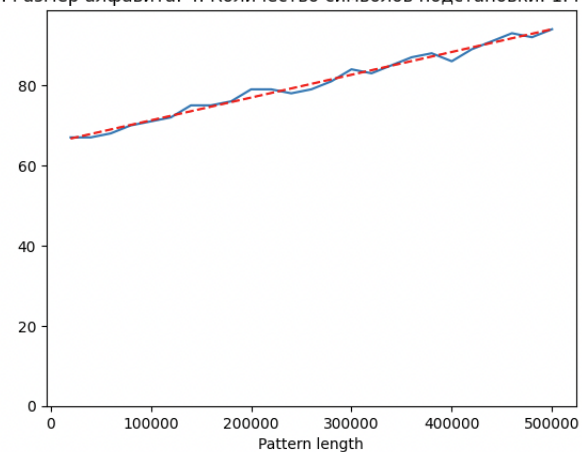
Размер текста: 1000000. Размер алфавита: 4. Количество символов подстановки: 0. КМП с уточненными границами



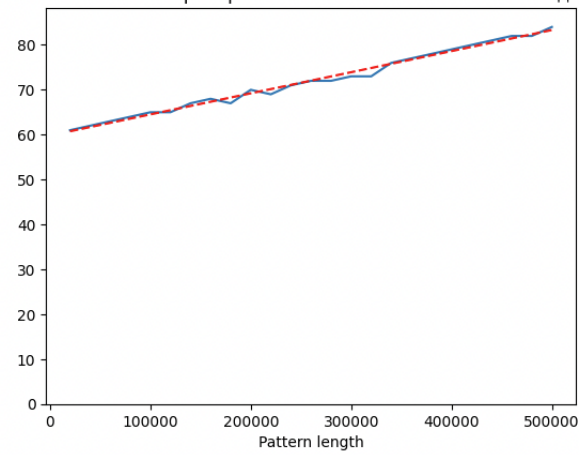
Размер текста: 1000000. Размер алфавита: 4. Количество символов подстановки: 1. КМП



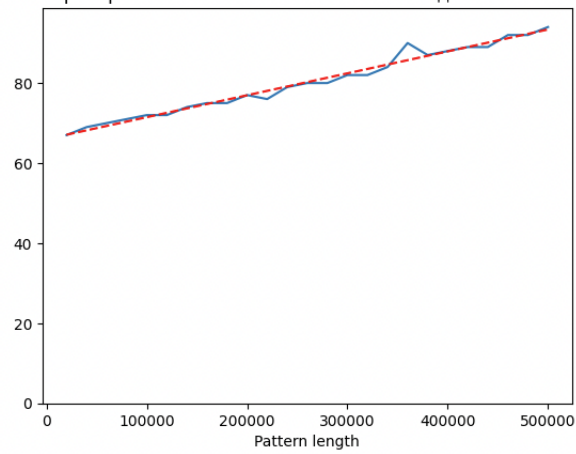
Размер текста: 1000000. Размер алфавита: 4. Количество символов подстановки: 1. КМП с уточненными границами



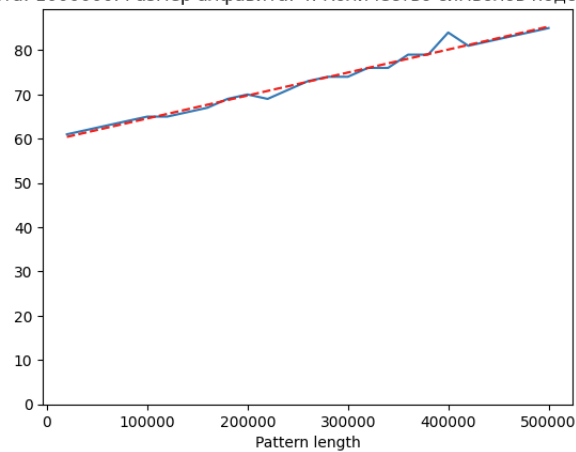
Размер текста: 1000000. Размер алфавита: 4. Количество символов подстановки: 2. КМП



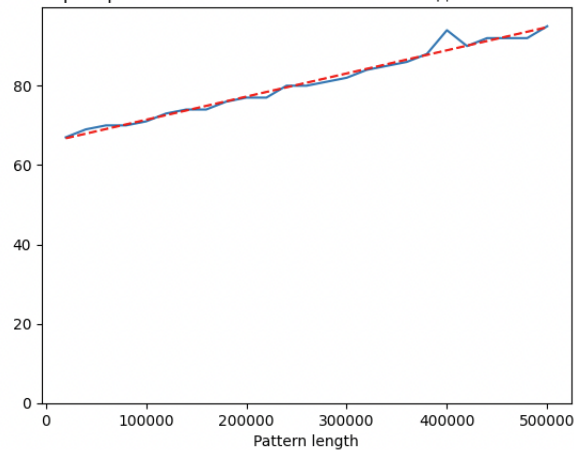
Размер текста: 1000000. Размер алфавита: 4. Количество символов подстановки: 2. КМП с уточненными границами



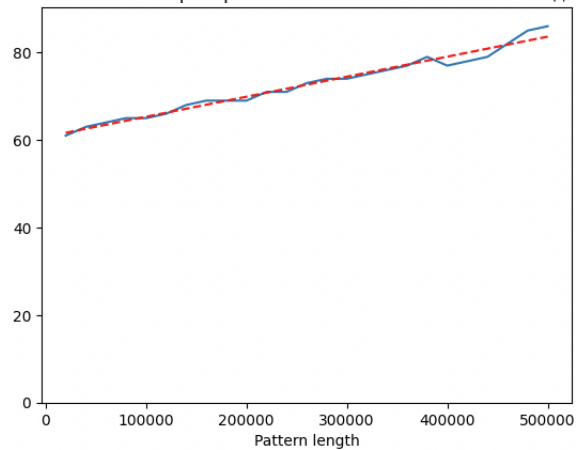
Размер текста: 1000000. Размер алфавита: 4. Количество символов подстановки: 3. КМП



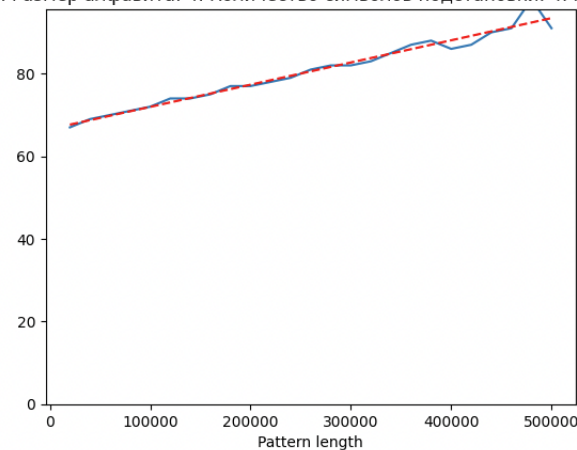
Размер текста: 1000000. Размер алфавита: 4. Количество символов подстановки: 3. КМП с уточненными границами



Размер текста: 1000000. Размер алфавита: 4. Количество символов подстановки: 4. КМП



Размер текста: 1000000. Размер алфавита: 4. Количество символов подстановки: 4. КМП с уточненными границами



На всех графиках видна довольно близкая к прямой картина (линии красным пуктиром - аппроксимация полиномом первой степени). Прямая имеет наклон, заметна линейная зависимость между размером паттерна и временем работы. Однако при минимальном размере паттерна сложность стремится к константе (не нулю) причем. Это позволяет говорить о том, что размер паттерна выступает слагаемым в первой степени в асимптотической формуле. Опять же,

использовав наши предположения о зависимости от размера текста, можем заключить, что сложность действительно:

$$O(text_size \cdot pattern_size)$$

4.3 Замечания

Оценка временной сложности алгоритма КМП с уточненными границами легко доказуема. С одной стороны, оценка сверху доказывается также, как и для обычного алгоритма КМП. С другой же, очевидно невозможно построить алгоритм, которому не пришлось бы как минимум считать входные данные. Поэтому сложность минимум $O(text_size + pattern_size)$.

Исследование на зависимость от размера алфавита не предполагались, судя по предлагаемому входным данным. Поэтому даже если это число влияет на сложность, оно предполагается константным коэффициентом, и, следовательно, опускается. Аналогично с количеством символов паттерна.

Все графики, различающиеся лишь количеством символов паттерна, почти идентичны. Можно сделать вывод, что количество символов паттерна, по крайней мере в таком небольшом количестве, практически не влияют на время работы алгоритма.

5. Выводы

- Алгоритмы были корректно реализованы
- Экспериментальный анализ временной сложности алгоритмов подтвердил теоретические оценки
- Влияние количества символов паттерна и размер алфавита невозможно исследовать с предлагаемыми входными данными.

6. Листинги

```
1 #include <iostream>
2 #include <vector>
3 #include <fstream>
4 #include <string>
5 #include <unordered_map>
6
7 inline bool compareSymbols(const char& i, const char& j) {
8     return i == j || i == '?' || j == '?';
9 }
10
11 std::vector<int> prefixFunction(const std::string& s) {
12     std::vector<int> pi(s.length());
13     pi[0] = 0;
14     int j;
15     for (int i = 1; i < static_cast<int>(s.length()); ++i) {
16         j = pi[i - 1];
17
18         while (j > 0 && !compareSymbols(s[i], s[j])) {
19             j = pi[j - 1];
20         }
21
22         if (compareSymbols(s[i], s[j])) {
23             j++;
24         }
25
26         pi[i] = j;
27     }
28     return pi;
29 }
30
31 std::vector<int> prefixFunctionSpecial(const std::string& s) {
32     std::vector<int> pi = prefixFunction(s);
33     for (int i = 0; i < static_cast<int>(pi.size()) - 1; ++i) {
34         if (pi[i + 1] == pi[i] + 1) {
35             pi[i] = 0;
36         }
37     }
38     return pi;
39 }
40
41 std::vector<int> kmp(const std::string& text, const std::string&
42     pattern, bool special = false) {
43     std::string data = pattern + "#" + text;
44     std::vector<int> pi = (special) ? prefixFunctionSpecial(data) :
45         prefixFunction(data);
46     std::vector<int> answers;
47     int pattern_size = static_cast<int>(pattern.length());
48     for (int i = 0; i < static_cast<int>(text.length()); ++i) {
49         if (pi[i + pattern_size + 1] == static_cast<int>(pattern.
50             length())) {
51             answers.push_back(i - pattern_size + 1);
52         }
53     }
54 }
```

```

51     return answers;
52 }
53
54 std::vector<int> naiveSearch(const std::string& text, const std::
55     string& pattern) {
56     std::vector<int> answers;
57     bool flag;
58     for (int i = 0; i < text.length(); ++i) {
59         flag = true;
60         for (int j = 0; j < pattern.length(); ++j) {
61             if (!compareSymbols(pattern[j], text[i + j])) {
62                 flag = false;
63             }
64             if (flag && j == pattern.length() - 1) {
65                 answers.push_back(i);
66             }
67         }
68     }
69     return answers;
70 }
71
72 std::string createRandomString(int size, int alphabet_size) {
73     std::string res;
74     for (int i = 0; i < size; ++i) {
75         auto t = rand() % alphabet_size;
76         res += 'A' + t;
77     }
78     return res;
79 }
80
81 struct Point {
82     int pattern_length;
83     long long time; // in milliseconds
84 };
85
86 void insertToFile(const std::vector<std::vector<Point>>& vv, const
87     std::string& file_name) {
88     std::ofstream file;
89     file.open(file_name);
90     std::string output;
91     for (const auto& v : vv) {
92         for (const auto& p : v) {
93             output += std::to_string(p.pattern_length) + ',' + std
94                 ::to_string(p.time) + ',';
95         }
96         output = output.substr(0, output.length() - 1);
97         output += '|';
98     }
99     output = output.substr(0, output.length() - 1);
100     file << output;
101     file.close();
102 }
103
104 long long process(std::function<void()> const& function) {
105     long long aggr = 0;

```

```

104     long long start, end;
105     for (int i = 0; i < 5; ++i) {
106         auto time_start = std::chrono::system_clock::now();
107         auto start_since_epoch = time_start.time_since_epoch();
108         start = std::chrono::duration_cast<std::chrono::
milliseconds>(start_since_epoch).count();
109         function();
110         auto time_end = std::chrono::system_clock::now();
111         auto end_since_epoch = time_end.time_since_epoch();
112         end = std::chrono::duration_cast<std::chrono::milliseconds
>(end_since_epoch).count();
113         aggr += (end - start);
114     }
115     return aggr / 5;
116 }
117
118 void doTest(int text_size, int alphabet_size, int
magic_symbols_amount = 0, bool with_naive = true, int
pattern_step = 100, int pattern_max_size = 3000, const std::
string& suffix = "") {
119     std::vector<Point> v_naive, v_kmp, v_kmp_special;
120     for (int pattern_size = pattern_step; pattern_size <=
pattern_max_size; pattern_size += pattern_step) {
121         std::string text = createRandomString(text_size,
alphabet_size);
122         std::string pattern = createRandomString(pattern_size,
alphabet_size);
123         if (magic_symbols_amount != 0) {
124             std::unordered_map<int, int> m;
125             std::vector<int> magic_indexes;
126             int max_size = pattern_size;
127             int randomized;
128             for (int i = 0; i < magic_symbols_amount; ++i) {
129                 randomized = rand() % (max_size - i);
130                 if (m.find(randomized) == m.end()) {
131                     magic_indexes.push_back(randomized);
132                     m[randomized] = max_size - 1;
133                 } else {
134                     magic_indexes.push_back(m[randomized]);
135                     m[randomized] = max_size - 1;
136                 }
137             }
138             for (const auto& i : magic_indexes) {
139                 pattern[i] = '?';
140             }
141         }
142
143         if (with_naive) {
144             auto naive_exec = [&text, &pattern]() { naiveSearch(
text, pattern); };
145             v_naive.push_back({pattern_size, process(naive_exec)});
146         }
147
148         auto kmp_exec = [&text, &pattern]() { kmp(text, pattern);};
149         v_kmp.push_back({pattern_size, process(kmp_exec)});
150

```



```

151     auto kmp_exec_special = [&text, &pattern]() { kmp(text,
152     pattern, true);};
153     v_kmp_special.push_back({pattern_size, process(
154     kmp_exec_special)});
155     }
156     std::vector<std::vector<Point>> vv = {v_kmp, v_kmp_special};
157     if (with_naive) {
158         vv.push_back(v_naive);
159     }
160     insertToFile(vv, "text_size_" + std::to_string(text_size) +
161     "_alphabet_size_"
162     + std::to_string(alphabet_size) +
163     "_magic_symbols_"
164     + std::to_string(
165     magic_symbols_amount) + "__" + suffix + ".txt");
166 }
167
168 int main() {
169     srand(time(NULL));
170
171     // 2 * 2 * 5 = 20 in total
172     // input data as it was requested
173     int count = 0;
174     int now;
175     int start = time(NULL);
176     for (const auto& text_size : {10000, 100000}) {
177         for (const auto& alphabet_size : {2, 4}) {
178             for (int i = 0; i < 5; ++i) {
179                 doTest(text_size, alphabet_size, i);
180                 count++;
181                 now = time(NULL);
182                 std::cout << (count * 5) << "% | " << ((now -
183     start) / 60) << " min from start\n";
184             }
185         }
186     }
187     std::cout << "requested data processed\n";
188
189     count = 0;
190     start = time(NULL);
191     // data is too small for kmp algorithms - need special tests
192     for (const auto& text_size : {1000000}) {
193         for (const auto& alphabet_size : {2, 4}) {
194             for (int i = 0; i < 5; ++i) {
195                 doTest(text_size, alphabet_size, i, false,
196     text_size / 50, text_size / 2, "kmp");
197                 count++;
198                 now = time(NULL);
199                 std::cout << (count * 10) << "% | " << ((now -
200     start) / 60) << " min from start\n";
201             }
202         }
203     }

```

Листинг 1: Тесты

```

1
2 file_names = []
3 text_sizes = [10000, 100000, 1000000]
4 alphabet_sizes = [2, 4]
5 prefixes = ["", "__kmp"]
6 magic_symbols = list(range(5))
7
8 for text_size in text_sizes:
9     for alphabet_size in alphabet_sizes:
10         for magic_symbol in magic_symbols:
11             for prefix in prefixes:
12                 file_names.append('text_size_' + str(text_size) +
13                                   '_alphabet_size_' + str(alphabet_size) +
14                                   '_magic_symbols_' + str(magic_symbol) +
15                                   prefix +
16                                   '.txt')
17
18 import os.path
19 graph_data = {}
20 for k, file_name in enumerate(file_names):
21     if os.path.isfile(file_name):
22         with open(file_name, 'r') as file:
23             print(file_name)
24             text = file.read()
25             lines = text.split('|')
26             lines = [[dot.split(',') for dot in line.split(';')] for line
27                     in lines]
28             graph_data[file_name] = lines
29
30 from matplotlib import pyplot as plt
31 import numpy as np
32
33
34 def extract_data_from_name(name):
35     paths = name.split('_')
36     return {
37         'text_size': paths[2],
38         'alphabet_size': paths[5],
39         'magic_symbols': paths[8] if len(paths) > 9 else paths[8][: -
40         len('.txt')],
41         'suffix': '' if len(paths) == 9 else ''
42     }
43
44 def create_title(data):
45     return f"Razmer teksta: {data['text_size']}. Razmer alphavita: {
46         data['alphabet_size']}. Kolichество simvolov podstanovki: {data
47         ['magic_symbols']}."

```

```

graph_data.items()):
48 XX = [[dot[0] for dot in line] for line in v]
49 YY = [[dot[1] for dot in line] for line in v]
50 plot_data = list(zip(XX, YY))
51 naive = plot_data[0]
52
53 title_data = extract_data_from_name(k)
54 ax = plt.subplot()
55 ax.plot(np.asarray(naive[0], int), np.asarray(naive[1], int))
56 ax.set_title(create_title(title_data) + " Naivniy algoritm.")
57 ax.set_xlabel('Pattern length')
58 plt.show()
59
60
61
62 counter = 0
63 kmp_names = ['KMP', 'KMP s utochnennimi granyami']
64 for k, v in filter(lambda data: len(data[0]) == len('
    text_size_1000000_alphabet_size_4_magic_symbols_4__kmp.txt'),
    graph_data.items()):
65
66 XX = [[dot[0] for dot in line] for line in v]
67 YY = [[dot[1] for dot in line] for line in v]
68 plot_data = list(zip(XX, YY))
69 for i, graph in enumerate(plot_data):
70     XX_asarr = np.asarray(graph[0], int)
71     YY_asarr = np.asarray(graph[1], int)
72
73     ax = plt.subplot()
74     ax.plot(XX_asarr, YY_asarr)
75
76     ax.set_xlabel('Pattern length')
77     plt.ylim([0, YY_asarr[-1] * 1.05])
78
79     title_data = extract_data_from_name(k)
80     ax.set_title(create_title(title_data) + ' ' + kmp_names[i])
81
82
83     p = np.polyfit(XX_asarr, YY_asarr, 1)
84
85     ya = np.polyval(p, XX_asarr)
86
87     plt.plot(XX_asarr, ya, '--', color='r')
88
89     plt.show()

```

Листинг 2: Построение графиков