

Raport z laboratorium 6

Filip Nikolow

27 maja 2021

1 Cel laboratorium

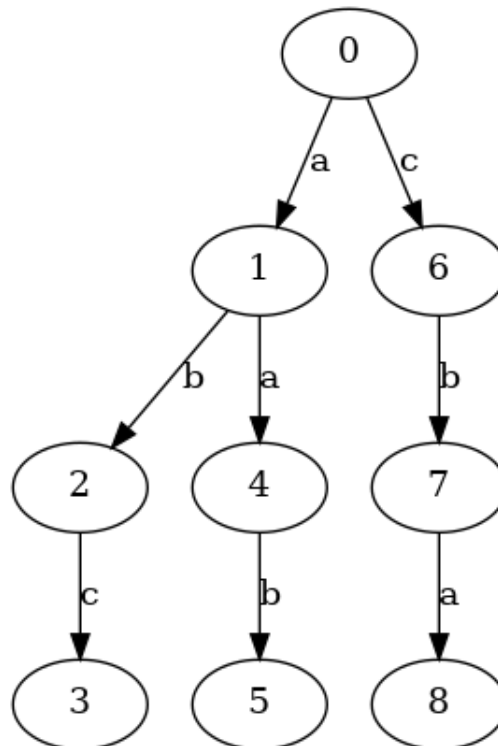
Celem laboratorium była implementacja i testy algorytmu wyszukiwania wzorca 2d.

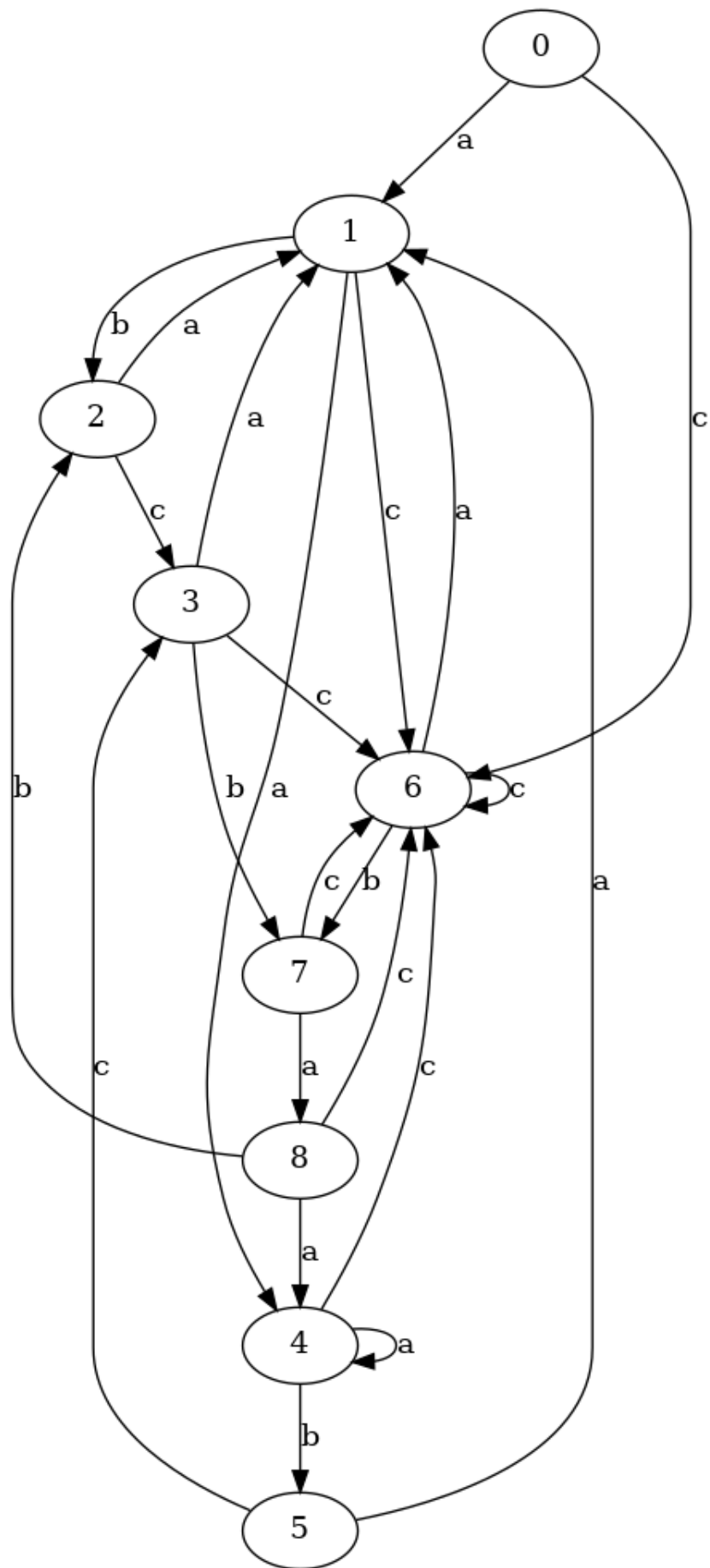
2 Realizacja poszczególnych poleceń

Kod do każdego z poleceń załączam na końcu sprawozdania.

2.1 Polecenie 1

Zaimplementowałem algorytm wyszukiwania wzorca 2d bazujący na automacie omawianym na ćwiczeniach. Poniżej załączam dwa zdjęcia: pierwsze to drzewo trie dla omawianego na zajęciach przykładowego wzorca, drugie to automat zbudowany na bazie tego drzewa (automat zupełny, bez pokazanej krawędzi domyślnej do stanu początkowego).





2.2 Polecenie 2

Użyłem zaimplementowanego algorytmu aby znaleźć wszystkie wystąpienia takich samych liter na takich samych pozycjach (za literę uznałem znak spełniający warunek `.isalpha()`). W poniższym wyniku działania programu każda krotka jest w formacie (linia, kolumna) i oznacza lewy górny róg wzorca:

```
'A': []
'B': []
'C': []
'D': []
'E': []
'F': []
'G': []
'I': []
'L': []
'M': []
'N': []
'O': []
'P': []
'Q': []
'R': []
'S': []
'T': []
'U': []
'V': []
'W': []
'X': []
'a': [(64, 2), (37, 4), (20, 6), (56, 11), (52, 12), (53, 12), (64,
14), (76, 21), (64, 22), (59, 24), (3, 30), (65, 35), (69, 35), (57,
36), (58, 36), (79, 37), (77, 42), (53, 48), (31, 50), (78, 59), (5,
60), (77, 61), (6, 63), (33, 66), (28, 69), (31, 73), (76, 74), (0,
82)]
'b': []
'c': [(41, 0), (68, 0), (13, 10), (82, 41), (10, 45), (3, 54)]
'd': [(37, 19)]
'e': [(10, 1), (14, 2), (24, 3), (17, 6), (76, 6), (77, 6), (80, 6),
(1, 8), (20, 10), (40, 11), (81, 14), (69, 15), (67, 17), (72, 23),
(40, 26), (18, 27), (73, 27), (51, 31), (42, 36), (29, 38), (71, 38),
(15, 43), (29, 43), (68, 46), (82, 47), (37, 48), (42, 48), (70, 49),
(47, 50), (58, 50), (46, 52), (22, 53), (57, 54), (58, 54), (41, 57),
(21, 61), (0, 63), (10, 64), (7, 65), (24, 65), (78, 65), (63, 66),
(28, 67), (65, 69), (66, 72), (28, 73), (59, 73), (4, 77)]
'f': [(77, 1), (30, 59)]
'g': []
'h': [(27, 2), (37, 2), (73, 12), (56, 31)]
'i': [(31, 0), (1, 5), (73, 13), (77, 13), (55, 17), (31, 31), (44,
33), (8, 37), (60, 45), (68, 51), (19, 55), (9, 60), (52, 69)]
'j': []
'k': []
'l': [(33, 45), (53, 45), (46, 61), (28, 72), (41, 77)]
'm': [(44, 0), (16, 5), (34, 40), (34, 60), (28, 70)]
'n': [(31, 1), (1, 9), (56, 13), (35, 18), (64, 29), (51, 32), (54,
```

```

33), (67, 35), (19, 37), (67, 40), (14, 54), (20, 56), (67, 57), (21,
62), (0, 83)]
'o': [(41, 1), (53, 1), (50, 2), (52, 8), (79, 10), (33, 11), (27,
17), (28, 17), (33, 26), (10, 27), (32, 34), (6, 38), (7, 38), (71,
42), (58, 45), (81, 52), (44, 55), (30, 58), (15, 60), (5, 66), (4,
75)]
'p': [(41, 18), (28, 71)]
'q': []
'r': [(1, 4), (52, 5), (33, 10), (7, 13), (17, 14), (15, 18), (69,
22), (43, 25), (67, 29), (60, 30), (33, 37), (47, 37), (6, 39), (62,
39), (55, 40), (46, 42), (6, 50), (19, 54), (20, 54), (28, 65), (31,
70)]
's': [(54, 0), (49, 14), (8, 21), (71, 24), (79, 24), (37, 34), (45,
34), (67, 37), (70, 41), (46, 44), (28, 45), (4, 49), (52, 53), (29,
56), (30, 56), (3, 57), (9, 58), (3, 63), (40, 63)]
't': [(37, 0), (50, 0), (16, 3), (71, 3), (72, 3), (23, 4), (24, 4),
(69, 5), (1, 6), (0, 7), (1, 7), (22, 8), (35, 10), (72, 10), (54,
11), (15, 12), (4, 14), (30, 16), (77, 22), (4, 23), (28, 23), (46,
24), (7, 29), (27, 31), (19, 33), (51, 33), (59, 33), (3, 37), (41,
45), (58, 49), (28, 52), (55, 54), (13, 55), (61, 56), (72, 59), (52,
61), (67, 71), (41, 73), (8, 75), (59, 75), (58, 78)]
'u': []
'v': []
'w': [(1, 3), (21, 70)]
'x': [(28, 68)]
'y': [(44, 5)]
'z': []

```

2.3 Polecenie 3

Załączam wynik działania, format jak w podpkt. 2:

```

th: []
t h: [(37, 0)]

```

2.4 Polecenie 4

Dokonałem wyszukiwania trzech wzorców:

r k o

Wyniki prezentuje jako haystack.png w którym znalezione wzorce są w negatywie, a dodatkowo wypisuje ilość znalezionych wzorców i ilość która powinna zostać znaleziona:

```

Found 298 occurrences of r.png
Should find 339
Found 21 occurrences of k.png
Should find 22
Found 310 occurrences of o.png
Should find 369

```

2.4.1 r

One of the simplest and natural types of information representation is by means of written texts. This type of data is characterized by the fact that it can be written down as a long sequence of characters. Such linear sequence is called a text. The texts are central in "word processing" systems, which provide facilities for the manipulation of texts. Such systems usually process objects that are quite large. For example, this book probably contains more than a million characters. Text algorithms occur in many areas of science and information processing. Many text editors and programming languages have facilities for processing texts. In biology, text algorithms arise in the study of molecular sequences. The complexity of text algorithms is also one of the central and most studied problems in theoretical computer science. It could be said that it is the domain in which practice and theory are very close to each other.

The basic textual problem in stringology is called pattern matching. It is used to access information and, no doubt, at this moment many computers are solving this problem as a frequently used operation in some application system. Pattern matching is comparable in this sense to sorting, or to basic arithmetic operations.

Consider the problem of a reader of the French dictionary "Garni Larousse," who wants all entries related to the name "Marie-Curie-Sklodowska." This is an example of a pattern matching problem, or string matching. In this case, the name "Marie-Curie-Sklodowska" is the pattern. Generally we may want to find a string called a pattern of length m inside a text of length n , where n is greater than m . The pattern can be described in a more complex way to denote a set of strings and not just a single word. In many cases n is very large. In genetics the pattern can correspond to a gene that can be very long; in image

The search of words or patterns in static texts is quite a different question than the previous pattern-matching mechanism. Dictionaries, for example, are organized in order to speed up the access to entries. Another example of the same question is given by indexes. Technical books often contain an index of chosen terms that gives pointers to parts of the text related to words in the index. The algorithms involved in the creation of an index form a specific group. The use of dictionaries or lexicons is often related to natural language processing. Lexicons of programming languages are small, and their representation is not a difficult problem during the development of a compiler. To the contrary, English contains approximately 100,000 words, and even twice that if inflected forms are considered. In French, inflected forms produce more than 700,000 words. The representation of lexicons of this size makes the problem a bit more challenging.

A simple use of dictionaries is illustrated by spelling checkers. The UNIX command, `spell`, reports the words in its input that are not stored in the lexicon. This rough approach does not yield a pertinent checker, but, practically, it helps to find typing errors. The lexicon used by `spell` contains approximately 70,000 entries stored within less than 60 kilobytes of random-access memory. Quick access to lexicons is a necessary condition for producing good parsers. The data structure useful for such access is called an index. In our book indexes correspond to data structures representing all factors of a given (presumably long) text. We consider problems related to the construction of such structures: suffix trees, directed acyclic word graphs, factor automata, suffix arrays. The PAT

tool developed at the N.O.E.D. Center (Waterloo, Canada) is an implementation of one of these structures tailored to work on large texts. There are several applications that effectively require some understanding of phrases in natural languages, such as data retrieval systems, interactive software, and character recognition.

An image scanner is a kind of photocopier. It is used to give a digitized version of an image. When the image is a page of text, the natural output of the scanner must be in a digital form available to a text editor. The transformation of a digitized image of a text into a usual computer representation of the text is realized by an Optical Character Recognition (OCR). Scanning a text with an OCR can be 50 times faster than retyping the text on a keyboard. Thus, OCR softwares are likely to become more common. But they still suffer from a high degree of imprecision. The average rate of error in the recognition of characters is approximately one percent. Even if this may happen to be rather small, this means that scanning a book produces approximately one error per line. This is compared with the usually very high quality of texts checked by specialists. Technical improvements on the hardware can help eliminate certain kinds of errors occurring on scanned texts in printed forms. But this cannot alleviate the problem associated with recognizing texts in printed forms. Reduction of the number of errors can thus only be achieved by considering the context of the characters, which assumes some understanding of the structure of the text. Image processing is related to the problem of two-dimensional pattern matching. Another related problem is the data structure for all subimages, which is discussed in this book in the context of the dictionary of basic factors.

The theoretical approach to the representation of lexicons is either by means of trees or finite state automata. It appears that both approaches are equally efficient. This shows the practical importance of the automata theoretic approach to text problems. At LITP (Paris) and IGM (Maison-la-Vallee) we have shown that the use of automata to represent lexicons is particularly efficient. Experiments have been done on a 700,000 word lexicon of LADL (Paris). The representation supports direct access to any word of the lexicon and takes only 300 kilobytes of random-access memory.

2.4.2 k

One of the simplest and natural types of information representation is by means of written texts. This type of data is characterized by the fact that it can be written down as a long sequence of characters. Such linear a sequence is called a text. The texts are central in "word processing" systems, which provide facilities for the manipulation of texts. Such systems usually process objects that are quite large. For example, this book probably contains more than a million characters. Text algorithms occur in many areas of science and information processing. Many text editors and programming languages have facilities for processing texts. In biology, text algorithms arise in the study of molecular sequences. The complexity of text algorithms is also one of the central and most studied problems in theoretical computer science. It could be said that it is the domain in which practice and theory are very close to each other.

The basic textual problem in stringology is called pattern matching. It is used to access information and, no doubt, at this moment many computers are solving this problem as a frequently used operation in some application system. Pattern matching is comparable in this sense to sorting, or to basic arithmetic operations.

Consider the problem of a reader of the French dictionary "Grand Larousse," who wants all entries related to the name "Marie-Curie-Skłodowska." This is an example of a pattern matching problem, or string matching. In this case, the name "Marie-Curie-Skłodowska" is the pattern. Generally we may want to find a string called a pattern of length m inside a text of length n , where n is greater than m . The pattern can be described in a more complex way to denote a set of strings and not just a single word. In many cases n is very large. In genetics the pattern can correspond to a gene that can be very long; in image

The search of words or patterns in static texts is quite a different question than the previous pattern-matching mechanism. Dictionaries, for example, are organized in order to speed up the access to entries. Another example of the same question is given by indexes. Technical books often contain an index of chosen terms that gives pointers to parts of the text related to words in the index. The algorithms involved in the creation of an index form a specific group. The use of dictionaries or lexicons is often related to natural language processing. Lexicons of programming languages are small, and their representation is not a difficult problem during the development of a compiler. To the contrary, English contains approximately 100,000 words, and even twice that if inflected forms are considered. In French, inflected forms produce more than 700,000 words. The representation of lexicons of this size makes the problem a bit more challenging.

A simple use of dictionaries is illustrated by spelling checkers. The UNIX command, `spell`, reports the words in its input that are not stored in the lexicon. This rough approach does not yield a pertinent checker, but, practically, it helps to find typing errors. The lexicon used by `spell` contains approximately 70,000 entries stored within less than 60 kilobytes of random-access memory. Quick access to lexicons is a necessary condition for producing good parsers. The data structure useful for such access is called an index. In our book indexes correspond to data structures representing all factors of a given (presumably long) text. We consider problems related to the construction of such structures: suffix trees, directed acyclic word graphs, factor automata, suffix arrays. The PAT

tool developed at the N.O.E.D. Center (Waterloo, Canada) is an implementation of one of these structures tailored to work on large texts. There are several applications that effectively require some understanding of phrases in natural languages, such as data retrieval systems, interactive software, and character recognition.

An image scanner is a kind of photocopier. It is used to give a digitized version of an image. When the image is a page of text, the natural output of the scanner must be in a digital form available to a text editor. The transformation of a digitized image of a text into a usual computer representation of the text is realized by an Optical Character Recognition (OCR). Scanning a text with an OCR can be 50 times faster than retyping the text on a keyboard. Thus, OCR softwares are likely to become more common. But they still suffer from a high degree of imprecision. The average rate of error in the recognition of characters is approximately one percent. Even if this may happen to be rather small, this means that scanning a book produces approximately one error per line. This is compared with the usually very high quality of texts checked by specialists. Technical improvements on the hardware can help eliminate certain kinds of errors occurring on scanned texts in printed forms. But this cannot alleviate the problem associated with recognizing texts in printed forms. Reduction of the number of errors can thus only be achieved by considering the context of the characters, which assumes some understanding of the structure of the text. Image processing is related to the problem of two-dimensional pattern matching. Another related problem is the data structure for all subimages, which is discussed in this book in the context of the dictionary of basic factors.

The theoretical approach to the representation of lexicons is either by means of trees or finite state automata. It appears that both approaches are equally efficient. This shows the practical importance of the automata theoretic approach to text problems. At LITP (Paris) and IGM (Marne-la-Vallée) we have shown that the use of automata to represent lexicons is particularly efficient. Experiments have been done on a 700,000 word lexicon of LADL (Paris). The representation supports direct access to any word of the lexicon and takes only 300 kilobytes of random-access memory.

2.4.3 o

One of the simplest and natural types of information representation is by means of written texts. This type of data is characterized by the fact that it can be written down as a long sequence of characters. Such linear a sequence is called a text. The texts are central in "word processing" systems, which provide facilities for the manipulation of texts. Such systems usually process objects that are quite large. For example, this book probably contains more than a million characters. Text algorithms occur in many areas of science and information processing. Many text editors and programming languages have facilities for processing texts. In biology, text algorithms arise in the study of molecular sequences. The complexity of text algorithms is also one of the central and most studied problems in the theoretical computer science. It could be said that it is the domain in which practice and theory are very close to each other.

The basic textual problem in stringology is called pattern matching. It is used to access information and, no doubt, at this moment many computers are solving this problem as a frequently used operation in some application system. Pattern matching is comparable in this sense to sorting, or to basic arithmetic operations.

Consider the problem of a reader of the French dictionary "Grand Larousse," who wants all entries related to the name "Marie-Curie-Skłodowska." This is an example of a pattern matching problem, or string matching. In this case, the name "Marie-Curie-Skłodowska" is the pattern. Generally we may want to find a string called a pattern of length m inside a text of length n , where n is greater than m . The pattern can be described in a more complex way to denote a set of strings and not just a single word. In many cases n is very large. In genetics the pattern can correspond to a gene that can be very long; in image

The search of words or patterns in static texts is quite a different question than the previous pattern-matching mechanism. Dictionaries, for example, are organized in order to speed up the access to entries. Another example of the same question is given by indexes. Technical books often contain an index of chosen terms that gives pointers to parts of the text related to words in the index. The algorithms involved in the creation of an index form a specific group. The use of dictionaries or lexicons is often related to natural language processing. Lexicons of programming languages are small, and their representation is not a difficult problem during the development of a compiler. To the contrary, English contains approximately 100,000 words, and even twice that if inflected forms are considered. In French, inflected forms produce more than 700,000 words. The representation of lexicons of this size makes the problem a bit more challenging.

A simple use of dictionaries is illustrated by spelling checkers. The UNIX command, `spell`, reports the words in its input that are not stored in the lexicon. This rough approach does not yield a pertinent checker, but, practically, it helps to find typing errors. The lexicon used by `spell` contains approximately 70,000 entries stored within less than 60 kilobytes of random-access memory. Quick access to lexicons is a necessary condition for producing good parsers. The data structure useful for such access is called an index. In our book indexes correspond to data structures representing all factors of a given (presumably long) text. We consider problems related to the construction of such structures: suffix trees, directed acyclic word graphs, factor automata, suffix arrays. The PAT

tool developed at the N O E D Center (Waterloo, Canada) is an implementation of one of these structures tailored to work on large texts. There are several applications that effectively require some understanding of phrases in natural languages, such as data retrieval systems, interactive software, and character recognition.

An image scanner is a kind of photocopier. It is used to give a digitized version of an image. When the image is a page of text, the natural output of the scanner must be in a digital form available to a text editor. The transformation of a digitized image of a text into a usual computer representation of the text is realized by an Optical Character Recognition (OCR). Scanning a text with an OCR can be 50 times faster than retyping the text on a keyboard. Thus, OCR softwares are likely to become more common. But they still suffer from a high degree of imprecision. The average rate of error in the recognition of characters is approximately one percent. Even if this may happen to be rather small, this means that a scanner produces approximately one error per line. This is compared with the usually very high quality of texts checked by specialists. Technical improvements on the hardware can help eliminate certain kinds of errors occurring on scanned texts in printed forms. But this cannot alleviate the problem associated with recognizing texts in printed forms. Reduction of the number of errors can thus only be achieved by considering the context of the characters, which assumes some understanding of the structure of the text. Image processing is related to the problem of two-dimensional pattern matching. Another related problem is the data structure for all subimages, which is discussed in this book in the context of the dictionary of basic factors.

The theoretical approach to the representation of lexicons is either by means of trees or finite state automata. It appears that both approaches are equally efficient. This shows the practical importance of the automata theoretic approach to text problems. At LITP (Paris) and IGM (Marne-la-Vallée) we have shown that the use of automata to represent lexicons is particularly efficient. Experiments have been done on a 700,000 word lexicon of LADL (Paris). The representation supports direct access to any word of the lexicon and takes only 300 kilobytes of random-access memory.

2.5 Polecenie 5

Podobnie jak w poleceniu 4, prezentuje wynik ze wzorcem w negatywie:

One of the simplest and natural types of information representation is by means of written texts. This type of data is characterized by the fact that it can be written down as a long sequence of characters. Such linear sequence is called a text. The texts are central in "word processing" systems, which provide facilities for the manipulation of texts. Such systems usually process objects that are quite large. For example, this book probably contains more than a million characters. Text algorithms occur in many areas of science and information processing. Many text editors and programming languages have facilities for processing texts. In biology, text algorithms arise in the study of molecular sequences. The complexity of text algorithms is also one of the central and most studied problems in theoretical computer science. It could be said that it is the domain in which practice and theory are very close to each other.

The basic textual problem in stringology is called pattern matching. It is used to access information and, no doubt, at this moment many computers are solving this problem as a frequently used operation in some application system. Pattern matching is comparable in this sense to sorting, or to basic arithmetic operations.

Consider the problem of a reader of the French dictionary "Grand Larousse," who wants all entries related to the name "Marie-Curie-Sklodowska." This is an example of a pattern matching problem, or string matching. In this case, the name "Marie-Curie-Sklodowska" is the pattern. Generally we may want to find a string called a pattern of length m inside a text of length n , where n is greater than m . The pattern can be described in a more complex way to denote a set of strings and not just a single word. In many cases n is very large. In genetics the pattern can correspond to a gene that can be very long; in image

The search of words or patterns in static texts is quite a different question than the previous pattern-matching mechanism. Dictionaries, for example, are organized in order to speed up the access to entries. Another example of the same question is given by indexes. Technical books often contain an index of chosen terms that gives pointers to parts of the text related to words in the index. The algorithms involved in the creation of an index form a specific group. The use of dictionaries or lexicons is often related to natural language processing. Lexicons of programming languages are small, and their representation is not a difficult problem during the development of a compiler. To the contrary, English contains approximately 100,000 words, and even twice that if inflected forms are considered. In French, inflected forms produce more than 700,000 words. The representation of lexicons of this size makes the problem a bit more challenging.

A simple use of dictionaries is illustrated by spelling checkers. The UNIX command, `spell`, reports the words in its input that are not stored in the lexicon. This rough approach does not yield a pertinent checker, but, practically, it helps to find typing errors. The lexicon used by `spell` contains approximately 70,000 entries stored within less than 60 kilobytes of random-access memory. Quick access to lexicons is a necessary condition for producing good parsers. The data structure useful for such access is called an index. In our book indexes correspond to data structures representing all factors of a given (presumably long) text. We consider problems related to the construction of such structures: suffix trees, directed acyclic word graphs, factor automata, suffix arrays. The PAT tool developed at the N.O.E.D. Center

(Waterloo, Canada) is an implementation of one of these structures tailored to work on large texts. There are several applications that effectively require some understanding of phrases in natural languages, such as data retrieval systems, interactive software, and character recognition.

An image scanner is a kind of photocopier. It is used to give a digitized version of an image. When the image is a page of text, the natural output of the scanner must be in a digital form available to a text editor. The transformation of a digitized image of a text into a usual computer representation of the text is realized by an Optical Character Recognition (OCR). Scanning a text with an OCR can be 50 times faster than retyping the text on a keyboard. Thus, OCR softwares are likely to become more common. But they still suffer from a high degree of imprecision. The average rate of error in the recognition of characters is approximately one percent. Even if this may happen to be rather small, this means that scanning a book produces approximately one error per line. This is compared with the usually very high quality of texts checked by specialists. Technical improvements on the hardware can help eliminate certain kinds of errors occurring on scanned texts in printed forms. But this cannot alleviate the problem associated with recognizing texts in printed forms. Reduction of the number of errors can thus only be achieved by considering the context of the characters, which assumes some understanding of the structure of the text. Image processing is related to the problem of two-dimensional pattern matching. Another related problem is the data structure for all subimages, which is discussed in this book in the context of the dictionary of basic factors.

The theoretical approach to the representation of lexicons is either by means of trees or finite state automata. It appears that both approaches are equally efficient. This shows the practical importance of the automata theoretic approach to text problems. At LITP (Paris) and IGM (Marne-la-Vallée) we have shown that the use of automata to represent lexicons is particularly efficient. Experiments have been done on a 700,000 word lexicon of LADL (Paris). The representation supports direct access to any word of the lexicon and takes only 300 kilobytes of random-access memory.

2.6 Polecenie 6

Poniżej załączam wzorce dla których przeprowadziłem pomiary, po kolei: small.png, medium.png, big.png.

One

Rys. 1: small.png

A simple use of dictionaries is illustrated by spelling checkers.

Rys. 2: medium.png

Consider the problem of a reader of the French dictionary "Grand Larousse," who wants all entries related to the name "Marie-Curie-Sklodowska." This is an example of a pattern matching problem, or string matching. In this case, the name "Marie-Curie-Sklodowska" is the pattern. Generally we may want to find a string called a pattern of length m inside a text of length n , where n is greater than m . The pattern can be described in a more complex way to denote a set of strings and not just a single word. In many cases n is very large. In genetics the pattern can correspond to a gene that can be very long; in image

Rys. 3: big.png

Poniżej załączam pomiary. Pomiary każdego przypadku przeprowadziłem 10-krotnie i wyciągnąłem z nich średnią:

```
=====Small pattern=====
Preprocessing time: 0.004592970500016236
Searching time: 1.0887323621998803
=====Medium pattern=====
Preprocessing time: 0.1500062629002059
Searching time: 0.9383168405998731
=====Big pattern=====
Preprocessing time: 2.218663644499975
Searching time: 0.9668237744997896
```

2.7 Polecenie 7

Załączam pomiary, podobnie jak w zadaniu 6, zostały powtórzone 10 razy, a w załączonych wynikach znajduje się średni czas wykonania.

```
=====Divided into 2 parts=====
Searching time: 0.9772503497999423
=====Divided into 4 parts=====
Searching time: 0.9691310410998994
=====Divided into 8 parts=====
Searching time: 0.9755429782999272
```

3 Kod

```
1  from PIL import Image
2  import numpy as np
3  import pydot
4  import tempfile
5  from queue import Queue
6  from timeit import default_timer as timer
7
8
9  class Node:
10     state_counter = 0
11
12     def __init__(self, parent=None):
13         self.state = Node.state_counter
14         Node.state_counter += 1
15
16         self.parent = parent
17         self.children = dict()
18         self.fail_link = parent
19
20     # finds next state, using fail links if necessary
21     def next_state(self, key):
22         if self.parent == self:
23             return self.children.get(key, self)
24         return self.children.get(key, self.fail_link.next_state(key))
25
26     def pretty_print(self, name=None, display=True):
27
28         def dfs_helper(node, graph, visited):
29             visited.add(node)
30             for k, v in node.children.items():
31                 graph.add_edge(pydot.Edge(str(node.state), str(v.state), label=str(k)))
32                 if v not in visited:
33                     dfs_helper(v, graph, visited)
34
35         graph = pydot.Dot(graph_type='digraph')
36         dfs_helper(self, graph, set())
37
38         if name is None:
39             fout = tempfile.NamedTemporaryFile(suffix=".png")
40             name = fout.name
41         else:
42             name += "_trie.png"
43         graph.write(name, format="png")
44         if display:
45             Image.open(name).show()
46
47
48  class PatternSearcher2d:
49     neutral = '#'
50
```

```

51     def __init__(self, pattern):
52         self.pattern = self.convert_text(pattern, self.neutral)
53         self.trie_root, self.accepting_states = None, None
54         self.text = None
55
56     def preprocess(self):
57         self.trie_root, self.accepting_states = self.aho_corasick()
58         return self
59
60     def load_text(self, text):
61         self.text = self.convert_text(text, self.neutral)
62         return self
63
64     # Neutral is a char non existent in the file we are going to scan
65     # Only important when supplying a list of lines
66     @staticmethod
67     def convert_text(text, neutral='#'):
68         if isinstance(text, list): # Assumes text is a list of lines
69             n = len(text)
70             m = max([len(line) for line in text])
71             arr = np.ndarray((n, m), dtype=object)
72             for i in range(n):
73                 for j in range(m):
74                     if j < len(text[i]):
75                         arr[i, j] = text[i][j]
76                     else:
77                         arr[i, j] = neutral
78             text = arr
79         elif not isinstance(text, (np.ndarray)):
80             Exception("Bad text format, has to be either list of lines or np.ndarray")
81         # Assumes text is a numpy color array of shape (n,m,color_data)
82         # where color_data is a list (so non_hashable)
83         else:
84             arr = np.ndarray(text.shape[:2], dtype=object)
85             for i in range(text.shape[0]):
86                 for j in range(text.shape[1]):
87                     arr[i, j] = tuple(text[i, j, :]) # 'Making' list hashable
88             text = arr
89         return text
90
91     # pattern must be of type np.ndarray
92     def aho_corasick(self, pattern=None):
93         if pattern is None:
94             pattern = self.pattern
95         (n, m) = pattern.shape[:2]
96         # Creating a simple trie
97         trie = Node()
98         trie.parent = trie
99         trie.fail_link = trie
100         accepting_states = []
101         alphabet = set()
102         for i in range(n):
103             p = trie

```

```

104         for j in range(m):
105             key = pattern[i, j]
106             if key not in p.children:
107                 alphabet.add(key)
108                 p.children[key] = Node(p)
109                 p = p.children[key]
110             accepting_states.append(p.state)
111         # trie.pretty_print("trie")
112         # Creating fail links
113         # Starting bfs from nodes with dist = 2, (dist={0,1} fail links are already correct)
114         Q = Queue()
115         [Q.put((k, v)) for c in trie.children.values() for k, v in c.children.items()]
116         while not Q.empty():
117             k, v = Q.get()
118             v.fail_link = v.parent.fail_link.next_state(k)
119             for k2, v2 in v.children.items():
120                 Q.put((k2, v2))
121
122         # Determinization of the automaton
123         Q = Queue()
124         Q.put(trie)
125         while not Q.empty():
126             v = Q.get()
127             for u in v.children.values():
128                 Q.put(u)
129             for k, u in v.fail_link.children.items():
130                 if k not in v.children:
131                     v.children[k] = u
132
133         # trie.pretty_print("automaton")
134         return trie, accepting_states
135
136     def search(self):
137         n, m = self.text.shape[:2]
138         pn, pm = self.pattern.shape[:2]
139         states = np.zeros(self.text.shape)
140         for i in range(n):
141             node = self.trie_root
142             for j in range(m):
143                 node = node.children.get(self.text[i, j], self.trie_root)
144                 states[i, j] = node.state
145         pattern_vert = self.convert_text([self.accepting_states])
146         root, accepting_state = self.aho_corasick(pattern_vert)
147         accepting_state = accepting_state[0]
148         found_coords = []
149         for j in range(m):
150             node = root
151             for i in range(n):
152                 node = node.children.get(states[i, j], root)
153                 if node.state == accepting_state:
154                     found_coords.append((i - pn + 1, j - pm + 1))
155         return found_coords
156

```

```

157     @staticmethod
158     def search_wrapper(text, pattern):
159         return PatternSearcher2d(pattern).preprocess().load_text(text).search()
160
161
162     def find_same_letter_pos(text):
163         alphabet = set()
164         for line in text:
165             for letter in line:
166                 alphabet.add(letter)
167         res = dict()
168         for letter in alphabet:
169             if letter.isalpha():
170                 res[letter] = PatternSearcher2d.search_wrapper(text, [letter, letter])
171         for k, v in sorted(res.items()):
172             print(repr(k) + ":", v)
173
174
175     def flip_colors(im, where, size=(20, 20)):
176         im = im.copy()
177         for x, y in where:
178             for i in range(size[0]):
179                 for j in range(size[1]):
180                     for k in range(3):
181                         im[x + i, y + j, k] = 255 - im[x + i, y + j, k]
182         return Image.fromarray(im.astype(np.uint8))
183
184
185     def bench(text, pattern, reps=10):
186         P = PatternSearcher2d(pattern)
187         P.load_text(text)
188
189         t1 = timer()
190         for _ in range(reps):
191             P.preprocess()
192         t2 = timer()
193         print("Preprocessing time:", (t2 - t1) / reps)
194
195         t1 = timer()
196         for _ in range(reps):
197             P.search()
198         t2 = timer()
199         print("Searching time:", (t2 - t1) / reps)
200
201
202     def divide_and_search(image, pattern, parts, reps=10):
203         P = PatternSearcher2d(pattern)
204         P.load_text(image)
205         P.preprocess()
206
207         size = image.shape[0]
208         parts = [
209             PatternSearcher2d.convert_text(image[i * size // parts:(i + 1) * size // parts, :])

```

```

210         for i in range(parts)
211     ]
212
213     t1 = timer()
214     for _ in range(reps):
215         for part in parts:
216             P.text = part
217             P.search()
218     t2 = timer()
219     print("Searching time:", (t2 - t1) / reps)
220
221
222 if __name__ == '__main__':
223     im = np.asarray(Image.open('../sources/haystack.png'))
224     # print(PatternSearcher2d.search_wrapper(['aaabcd', 'eeaab', 'ppcba'],
225     #                                         ["abc", "aab", "cba"]))
226     with open("../sources/haystack.txt", "r") as f:
227         text = f.readlines()
228         # Task2
229         find_same_letter_pos(text)
230         # Task3
231         print("th:", PatternSearcher2d.search_wrapper(text, ['th', 'th']))
232         print("t h:", PatternSearcher2d.search_wrapper(text, ['t h', 't h']))
233         # Task4
234         r = np.asarray(Image.open('../patterns/r.png'))
235         k = np.asarray(Image.open('../patterns/k.png'))
236         o = np.asarray(Image.open('../patterns/o.png'))
237         where = PatternSearcher2d.search_wrapper(im, r)
238         flip_colors(im, where, r.shape).show()
239         print("Found", len(where), "occurences of r.png")
240         print("Should find", len(PatternSearcher2d.search_wrapper(text, ["r"])))
241         where = PatternSearcher2d.search_wrapper(im, k)
242         flip_colors(im, where, k.shape).show()
243         print("Found", len(where), "occurences of k.png")
244         print("Should find", len(PatternSearcher2d.search_wrapper(text, ["k"])))
245         where = PatternSearcher2d.search_wrapper(im, o)
246         flip_colors(im, where, o.shape).show()
247         print("Found", len(where), "occurences of o.png")
248         print("Should find", len(PatternSearcher2d.search_wrapper(text, ["o"])))
249         # Task5
250         pattern = np.asarray(Image.open('../patterns/pattern.png'))
251         where = PatternSearcher2d.search_wrapper(im, pattern)
252         flip_colors(np.asarray(Image.open('../sources/haystack.png')),
253                     where,
254                     size=pattern.shape).show()
255         # Task6
256         print("=====Small pattern=====")
257         bench(im, np.asarray(Image.open("../patterns/small.png")), 10)
258         print("=====Medium pattern=====")
259         bench(im, np.asarray(Image.open("../patterns/medium.png")), 10)
260         print("=====Big pattern=====")
261         bench(im, np.asarray(Image.open("../patterns/big.png")), 10)
262         # Task7

```

```
263     for s in [2, 4, 8]:
264         print("====Diveded into", s, "parts====")
265         divide_and_search(im, np.asarray(Image.open("../patterns/medium.png")), s)
```
