

# Project 3

Mengxiang Jiang  
CSEN 5336 Analysis of Algorithms

October 26, 2023

**Problem 1.** Five DNA sequences of different lengths are given (DNA1.txt, DNA2.txt, etc).

Find the occurrences of the pattern ACT, GATTACA, and CATCATCAT in each using the naive algorithm. Submit your code to solve the problem. How much time (in seconds or milliseconds) is required by your computer to run the algorithm?

```
1 def naive(text, pattern):
2     matches = []
3     for i in range(0, len(text) - len(pattern)):
4         found = True
5         for j in range(len(pattern)):
6             if text[i+j] != pattern[j]:
7                 found = False
8                 break
9         if found == True:
10             matches.append(i)
11     return matches
```

Listing 1: naive algorithm

```
1 import time
2 if __name__ == '__main__':
3     text = ''
4     patterns = ['ACT', 'GATTACA', 'CATCATCATCAT']
5     for fname in ['DNA1.txt', 'DNA2.txt', 'DNA3.txt', 'DNA4.txt', 'DNA5.txt']:
6         with open(fname, 'r') as f:
7             text = f.read().rstrip()
8         for pattern in patterns:
9             start_time = time.perf_counter()
10            matches = naive(text, pattern)
11            end_time = time.perf_counter()
12            print(f'For {fname} and pattern {pattern}, the naive algorithm
13                took {end_time-start_time} seconds with {len(matches)} matches')
```

Listing 2: main function code

Note: the main function code is similar for running all the algorithms.

```

1 $ python naive.py
2   For DNA1.txt and pattern ACT, the naive algorithm took
   0.3632765369984554 seconds with 65495 matches
3   For DNA1.txt and pattern GATTACA, the naive algorithm took
   0.3577618750023248 seconds with 242 matches
4   For DNA1.txt and pattern CATCATCATCAT, the naive algorithm took
   0.3621113970002625 seconds with 0 matches
5   For DNA2.txt and pattern ACT, the naive algorithm took
   0.727408620998176 seconds with 131034 matches
6   For DNA2.txt and pattern GATTACA, the naive algorithm took
   0.7238480179985345 seconds with 491 matches
7   For DNA2.txt and pattern CATCATCATCAT, the naive algorithm took
   0.7162508259971219 seconds with 1 matches
8   For DNA3.txt and pattern ACT, the naive algorithm took
   1.4487447140018048 seconds with 262192 matches
9   For DNA3.txt and pattern GATTACA, the naive algorithm took
   1.4347486420010682 seconds with 1002 matches
10  For DNA3.txt and pattern CATCATCATCAT, the naive algorithm took
   1.4541748980009288 seconds with 1 matches
11  For DNA4.txt and pattern ACT, the naive algorithm took
   2.882448870997905 seconds with 524454 matches
12  For DNA4.txt and pattern GATTACA, the naive algorithm took
   2.870233528999961 seconds with 1976 matches
13  For DNA4.txt and pattern CATCATCATCAT, the naive algorithm took
   2.910066188000201 seconds with 5 matches
14  For DNA5.txt and pattern ACT, the naive algorithm took
   5.771988509000948 seconds with 1050689 matches
15  For DNA5.txt and pattern GATTACA, the naive algorithm took
   5.708121202998882 seconds with 4161 matches
16  For DNA5.txt and pattern CATCATCATCAT, the naive algorithm took
   5.878649592999864 seconds with 3 matches

```

Listing 3: Terminal output running code

Even though the time complexity of the naive algorithm is  $O((n-m+1)m)$  in the worst case, in practice the probability the algorithm will break out early makes it perform significantly faster than expected where the times seem to indicate something closer to  $O(n)$ .

**Problem 2.** Repeat the previous problem with the Rabin-Karp algorithm.

```
1 SigMap = {'A': 0, 'T': 1, 'G': 2, 'C': 3}
2 PriMod = 13
3
4 def rabin_karp(text, pattern):
5     matches = []
6     plen = len(pattern)
7     phash = rabin_hash(pattern)
8     for i in range(0, len(text) - len(pattern)):
9         if i == 0:
10             thash = rabin_hash(text[0:len(pattern)])
11         else:
12             thash = roll(plen, thash,
13                          text[i - 1], text[i + plen - 1])
14         if thash == phash:
15             if text[i:i+plen] == pattern:
16                 matches.append(i)
17     return matches
18
19 def rabin_hash(text):
20     rhash = 0
21     for c in text:
22         rhash = (rhash * len(SigMap) % PriMod + SigMap[c]) % PriMod
23     return rhash
24
25 def roll(plen, oldhash, oldchar, newchar):
26     newhash = ((oldhash + PriMod
27                - SigMap[oldchar] * (len(SigMap) ** (plen - 1)) % PriMod)
28                * len(SigMap) + SigMap[newchar]) % PriMod
29     return newhash
```

Listing 4: Rabin-Karp algorithm

```

1 $ python rabin_karp.py
2   For DNA1.txt and pattern ACT, the Rabin-Karp algorithm took
   0.639081138993788 seconds with 65495 matches
3   For DNA1.txt and pattern GATTACA, the Rabin-Karp algorithm took
   0.6668751739998697 seconds with 242 matches
4   For DNA1.txt and pattern CATCATCATCAT, the Rabin-Karp algorithm took
   0.7319347989978269 seconds with 0 matches
5   For DNA2.txt and pattern ACT, the Rabin-Karp algorithm took
   1.2784735669993097 seconds with 131034 matches
6   For DNA2.txt and pattern GATTACA, the Rabin-Karp algorithm took
   1.3436157860051026 seconds with 491 matches
7   For DNA2.txt and pattern CATCATCATCAT, the Rabin-Karp algorithm took
   1.4682704229999217 seconds with 1 matches
8   For DNA3.txt and pattern ACT, the Rabin-Karp algorithm took
   2.5847115970027517 seconds with 262192 matches
9   For DNA3.txt and pattern GATTACA, the Rabin-Karp algorithm took
   2.6529827709964593 seconds with 1002 matches
10  For DNA3.txt and pattern CATCATCATCAT, the Rabin-Karp algorithm took
   2.936001311005384 seconds with 1 matches
11  For DNA4.txt and pattern ACT, the Rabin-Karp algorithm took
   5.080261541996151 seconds with 524454 matches
12  For DNA4.txt and pattern GATTACA, the Rabin-Karp algorithm took
   5.316721924995363 seconds with 1976 matches
13  For DNA4.txt and pattern CATCATCATCAT, the Rabin-Karp algorithm took
   5.90457762600272 seconds with 5 matches
14  For DNA5.txt and pattern ACT, the Rabin-Karp algorithm took
   10.213705198999378 seconds with 1050689 matches
15  For DNA5.txt and pattern GATTACA, the Rabin-Karp algorithm took
   10.634073597997485 seconds with 4161 matches
16  For DNA5.txt and pattern CATCATCATCAT, the Rabin-Karp algorithm took
   11.496143873002438 seconds with 3 matches

```

Listing 5: Terminal output running code

Despite having better expected time complexity of  $O(n)$  compared to the naive algorithm, Rabin-Karp appears to be about twice as slow for these examples.

**Problem 3.** Repeat the previous problem with the finite automata algorithm.

```
1 SigMap = {'A': 0, 'T': 1, 'G': 2, 'C': 3}
2
3 def finite_automata(text, pattern):
4     matches = []
5     m = len(pattern)
6     tf = computeTF(pattern)
7
8     state = 0
9     for i in range(len(text)):
10         state = tf[state][SigMap[text[i]]]
11         if state == m:
12             matches.append(i-m+1)
13
14     return matches
15
16 def computeTF(pattern):
17     m = len(pattern)
18     tf = [[0 for i in range(len(SigMap))] for j in range(m + 1)]
19     for state in range(m + 1):
20         for c in SigMap:
21             nextState = getNextState(pattern, state, c)
22             tf[state][SigMap[c]] = nextState
23     return tf
24
25 def getNextState(pattern, state, c):
26     # c matches next character in pattern
27     if state < len(pattern) and c == pattern[state]:
28         return state + 1
29
30     # find longest prefix which is also the suffix
31     # of pattern[0:nextState]
32     i = 0
33     for nextState in range(state, 0, -1):
34         if pattern[nextState-1] == c:
35             while (i < nextState - 1):
36                 pfxidx = state - nextState + 1 + i
37                 if pattern[i] != pattern[pfxidx]:
38                     break
39                 i += 1
40             if i == nextState - 1:
41                 return nextState
42     return 0
```

Listing 6: finite automata algorithm

```

1 $ python finite_automata.py
2   For DNA1.txt and pattern ACT, the finite automata algorithm took
   0.15551975400012452 seconds with 65495 matches
3   For DNA1.txt and pattern GATTACA, the finite automata algorithm took
   0.14494252900476567 seconds with 242 matches
4   For DNA1.txt and pattern CATCATCATCAT, the finite automata algorithm
   took 0.14850315799412783 seconds with 0 matches
5   For DNA2.txt and pattern ACT, the finite automata algorithm took
   0.30578777501068544 seconds with 131034 matches
6   For DNA2.txt and pattern GATTACA, the finite automata algorithm took
   0.29742643699864857 seconds with 491 matches
7   For DNA2.txt and pattern CATCATCATCAT, the finite automata algorithm
   took 0.2898057919956045 seconds with 1 matches
8   For DNA3.txt and pattern ACT, the finite automata algorithm took
   0.6287656860076822 seconds with 262192 matches
9   For DNA3.txt and pattern GATTACA, the finite automata algorithm took
   0.59431980199588 seconds with 1002 matches
10  For DNA3.txt and pattern CATCATCATCAT, the finite automata algorithm
   took 0.5852815399994142 seconds with 1 matches
11  For DNA4.txt and pattern ACT, the finite automata algorithm took
   1.2557424819970038 seconds with 524454 matches
12  For DNA4.txt and pattern GATTACA, the finite automata algorithm took
   1.202983769995626 seconds with 1976 matches
13  For DNA4.txt and pattern CATCATCATCAT, the finite automata algorithm
   took 1.1666112150123809 seconds with 5 matches
14  For DNA5.txt and pattern ACT, the finite automata algorithm took
   2.4493759730103193 seconds with 1050689 matches
15  For DNA5.txt and pattern GATTACA, the finite automata algorithm took
   2.3914340260089375 seconds with 4161 matches
16  For DNA5.txt and pattern CATCATCATCAT, the finite automata algorithm
   took 2.339489956997568 seconds with 3 matches

```

Listing 7: Terminal output running code

The finite automata algorithm performs faster than both the naive and Rabin-Karp algorithms, and also seems truly independent of the length of the pattern, performing about the same or even quicker in the longer patterns.

**Problem 4.** Repeat the previous problem with the Knuth-Morris-Pratt(KMP) algorithm.

```
1 def kmp(text, pattern):
2     matches = []
3     m = len(pattern)
4     n = len(text)
5     pfxfun = getPrefixFunction(pattern)
6     i = 0
7     j = 0
8
9     while i < n:
10         if text[i] == pattern[j]:
11             i += 1
12             j += 1
13             if j == m:
14                 matches.append(i-j)
15                 j = pfxfun[j-1]
16         else:
17             if j > 0:
18                 j = pfxfun[j-1]
19             else:
20                 i += 1
21
22     return matches
23
24 def getPrefixFunction(pattern):
25     pfxfun = [0 for i in range(len(pattern))]
26     i = 1
27     j = 0
28     while i < len(pattern):
29         if pattern[i] == pattern[j]:
30             j += 1
31             pfxfun[i] = j
32             i += 1
33         else:
34             if j != 0:
35                 j = pfxfun[j - 1]
36             else:
37                 pfxfun[i] = 0
38                 i += 1
39     return pfxfun
```

Listing 8: Knuth-Morris-Pratt algorithm

```

1  $ python kmp.py
2  For DNA1.txt and pattern ACT, the Knuth-Morris-Pratt algorithm took
   0.2486665429896675 seconds with 65495 matches
3  For DNA1.txt and pattern GATTACA, the Knuth-Morris-Pratt algorithm
   took 0.2490224240027601 seconds with 242 matches
4  For DNA1.txt and pattern CATCATCATCAT, the Knuth-Morris-Pratt
   algorithm took 0.24674882799445186 seconds with 0 matches
5  For DNA2.txt and pattern ACT, the Knuth-Morris-Pratt algorithm took
   0.4929636379965814 seconds with 131034 matches
6  For DNA2.txt and pattern GATTACA, the Knuth-Morris-Pratt algorithm
   took 0.4923475959949428 seconds with 491 matches
7  For DNA2.txt and pattern CATCATCATCAT, the Knuth-Morris-Pratt
   algorithm took 0.49851692000811454 seconds with 1 matches
8  For DNA3.txt and pattern ACT, the Knuth-Morris-Pratt algorithm took
   0.9728445020009531 seconds with 262192 matches
9  For DNA3.txt and pattern GATTACA, the Knuth-Morris-Pratt algorithm
   took 0.7028228560084244 seconds with 1002 matches
10 For DNA3.txt and pattern CATCATCATCAT, the Knuth-Morris-Pratt
   algorithm took 0.704158250009641 seconds with 1 matches
11 For DNA4.txt and pattern ACT, the Knuth-Morris-Pratt algorithm took
   1.40684947000409 seconds with 524454 matches
12 For DNA4.txt and pattern GATTACA, the Knuth-Morris-Pratt algorithm
   took 1.424388044004445 seconds with 1976 matches
13 For DNA4.txt and pattern CATCATCATCAT, the Knuth-Morris-Pratt
   algorithm took 1.4213517800089903 seconds with 5 matches
14 For DNA5.txt and pattern ACT, the Knuth-Morris-Pratt algorithm took
   2.793876569994609 seconds with 1050689 matches
15 For DNA5.txt and pattern GATTACA, the Knuth-Morris-Pratt algorithm
   took 2.8016387920069974 seconds with 4161 matches
16 For DNA5.txt and pattern CATCATCATCAT, the Knuth-Morris-Pratt
   algorithm took 2.790721905999817 seconds with 3 matches

```

Listing 9: Terminal output running code

The Knuth-Morris-Pratt algorithm performs faster than both the naive and Rabin-Karp algorithms but slower than the finite automata. It also seems independent of the length of the pattern like finite automata. It might be slower than finite automata because it does more comparisons per character of the text.