

CÔTE D'AZUR	Pattern Matching Exact	CÔTE D'AZUR	Rabin Karp algorithm
Algorithmics for Biology Jean-Paul Complexity Pat. Matching Rabin-Karp Antomata KMP BM Graphs Dyn.Prog. Sequences	Exact search for a word of length m in a text of length n ≫ m. compare the m letters of the word with the m letters of the text beginning at position j, for j = 1,n - m + 1. complexity : O(m × n) 1 n = long[T] m = long[P] for (s=0, s<n-m; <="" li="" s++):=""> if P[0m-1] = T[ss+m-1] then print ``the pattern is present at position'', s; </n-m;> To improve the algorithm, information from step j or from previous steps must be taken into account at step j + 1.	Algorithmics for Biology Jean-Paul Complexity Pat. Matching Rabin-Karp Antomata KMP BM Graphs Dyn.Prog. Sequences	Rabin Karp algorithm = integer encoding of substrings• Worst-case execution time : $O((n - m + 1)m)$ • Average execution time good.Alphabet : Σ Word on Σ : string of k consecutive characters $d = \Sigma $ Word on Σ : number written in base d of length k.Here : $\Sigma = \{0, 1, 2,9\}$ Pattern P[1m] : we note p its corresponding decimal value.Text T[1n] : we note t_s the decimal value of the substring T[s+1s+m]• Computation of p in $O(m)$ Horner's scheme : $p = P[m] + 10(P[m-1] + 10(P[m-2] + + 10(P[2] + 10P[1])))$
	 Example : Let us consider the pattern P = ATAAG If the pattern is present in position i, then we can deduce that in position i + 1, i + 2, i + 3 and i + 4 the pattern cannot appear. If the letter of the text at i is an A, but the pattern is not present in position i, then the pattern cannot be present in position i + 1 (perhaps in position i + 2). 		• Computation of t_1 in $O(m)$ • Computation of t_{s+1} : $t_{s+1} = 10(t_s - 10^{m-1}T[s]) + T[s+m]$. Example : text $\equiv 134512$ $m = 5$ $t_1 = 13451$ $t_2 = (13451 - 10^4 \times 1) \times 10 + 2 = 34512$. Computation of constant 10^{m-1} in $O(log_2(m))$ $\implies t_0, t_1,, t_{n-m}$ can be theoretically computed in $O(n+m)$ \implies occurrences of P in $T[1n]$ can be computed in $O(n+m)$
	 < 마 > 4 문 > 4 문 > 4 문 > 로 - 카익은 20/112 		\implies occurrences of P in T[1n] can be computed in $O(n+m)$

CÔTE D'AZUR Rabin Karp algorithm : Problems

Problem #1

Algorithmics for Biology

Rabin-Karp

Algorithmics for Biology

If one of the numbers $t_0, t_1, ..., t_{n-m}$ is too large, arithmetic operations on *m* digits no longer take a constant time.

Remedy : computing modulo q.

 $\overline{p, t_0, t_1} \dots t_{n-m}$ modulo q can be computed in O(n+m). We choose q such that $10 \times q$ (in fact $d \times q$) just fits on a machine word. The formula for calculating t_{s+1} becomes :

 $t_{s+1} = (10(t_s - T[s]10^{m-1} \mod q) + T[s+m]) \mod q$

Problem #2

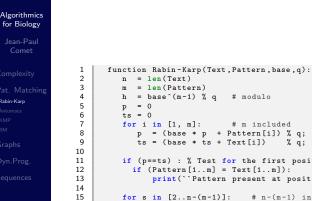
Even if $t_s = p$ implies the equality $(t_s \mod q) = (p \mod q)$, the fact of having calculated the values of t_s modulo q, makes the test insufficient :

• $(t_s \mod q) = (p \mod q)$ does not imply $t_s = p$

Remedy : We then use $(t_s \mod q) = (p \mod q)$ as a quick test, and when the moduli are equal, we test each letter.

◆□▶ < @ ▶ < E ▶ < E ▶ ○ 22/112</p>

CÔTE D'AZUR Rabin Karp algorithm : pseudocode



16

17 18

19

for Biology

Rabin-Karp

1	h = base^(m-1) % q # modulo			
5	p = 0			
5	ts = 0			
	for i in [1, m]: # m included			
3	p = (base * p + Pattern[i]) % q;			
)	ts = (base * ts + Text[i]) % q;			
)				
L	if (p==ts) : % Test for the first position			
2	if (Pattern[1m] = Text[1m]):			
3	<pre>print(``Pattern present at position '', 1);</pre>			
1				
5	<pre>for s in [2n-(m-1)]: # n-(m-1) included</pre>			
5	ts = (base*(ts - Text[s]*h)+ Text[s+m]) % q;			
7	<pre>if (p==ts):</pre>			
2	if (Pattern[1 m] = Text[s s+m-1]).			

(Pattern[1..m] = Text[s..s+m-1]): print(``Pattern present at position '', s)

<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

CÔTE D'AZUR Rabin Karp algorithm : complexity

Worst-case complexity calculation :

- Choose an example where you spend the whole time comparing letter to letter.
- Example : $P = a^m$ and $T = a^n$
- $\Rightarrow \Theta((n-m+1) \times m).$

(quadratic)

Estimated chance of having $t_s = p \mod q : \frac{1}{q}$. In fact, we have one chance in q of choosing $p \mod q$.

Average complexity calculation :

- Lines 2-13 (computation of p,t_1) : 0(m)
- Lines 15-19 (without any letter-to-letter test) : O(n)
- Lines 17-19 (letter-to-letter tests) : $O(m(v + \frac{n}{a})) = O(mv + \frac{n}{a} \times m)$ where v is the number of occurrences of the pattern.
- \Rightarrow Average-case complexity : $0(m + n + m(v + \frac{n}{a}))$

If v is small (O(1) i.e. \sim one occurrence of the pattern) and if q > m: Average-case complexity : O(m + n).