

1. Ponazorite delovanje KMP algoritma tako, da poiščete vzorec $p = ACACA$ v besedilu $t = ACACACGACACACA$.

2	0	1	2	3	4	2	0	0	1	2	0	0	1	2	3	4	5	3	4	5
i	1	2	3	4	5	5	5	6	7	8	8	9	10	11	12	13	14	14	15	
p[i+1]	A	C	A	C	A	A	A	A	C	A	A	A	C	A	C	A	E	C	A	
t[i]	A	C	A	C	T	T	T	A	C	G	G	A	C	A	C	A	C	C	A	
match	✓	✓	✓	✓	✗	✗	✗	✓	✓	✗	✗	✓	✓	✓	✓	✓	✓	✓	✓	

p	A	C	A	C	A
π	0	0	1	2	3
	1	2	3	4	5

$q=0$ in znaka zaporedij nista enaka (while loop): vrni za $q = \pi[q]$

$\pi[4]$ je 2. Spet preveri pogoj.

$q=0$ zato se while preneha. Znaka nista enaka, zato gre v naslednjo iteracijo for zanke \Rightarrow poveča i , a ne tudi q .

def KMP(p, t):

$n = |t|, m = |p|, r = []$

$\pi = \text{constructPi}(p)$

$q = 0$

for $i = 1 \dots n$ do

while $q > 0$ and $p[q+1] \neq t[i]$ do

$q = \pi[q]$

if $p[q+1] = t[i]$ then

$q++$

if $q = m$ then

$r += r[i-q+1]$

return r

$q=m$ zato v r doda indeks znaka, ki je začel ujemanje (t.j. v tem primeru $i-q+1 = 15-5+1 = 11$)

zakaj ne 14? Takoj ko se poveča q se to preveri, preden se poveča i .

2. Zgradite priponsko drevo za tekst $t = ATTTGCCG$.

- (a) Ali tekst vsebuje vzorec $p = TT$? Kje?
(b) Ali tekst vsebuje vzorec $p = CGT$? Kje?

1. Izpišeš vse pripone besedila: $(|t|+1)$

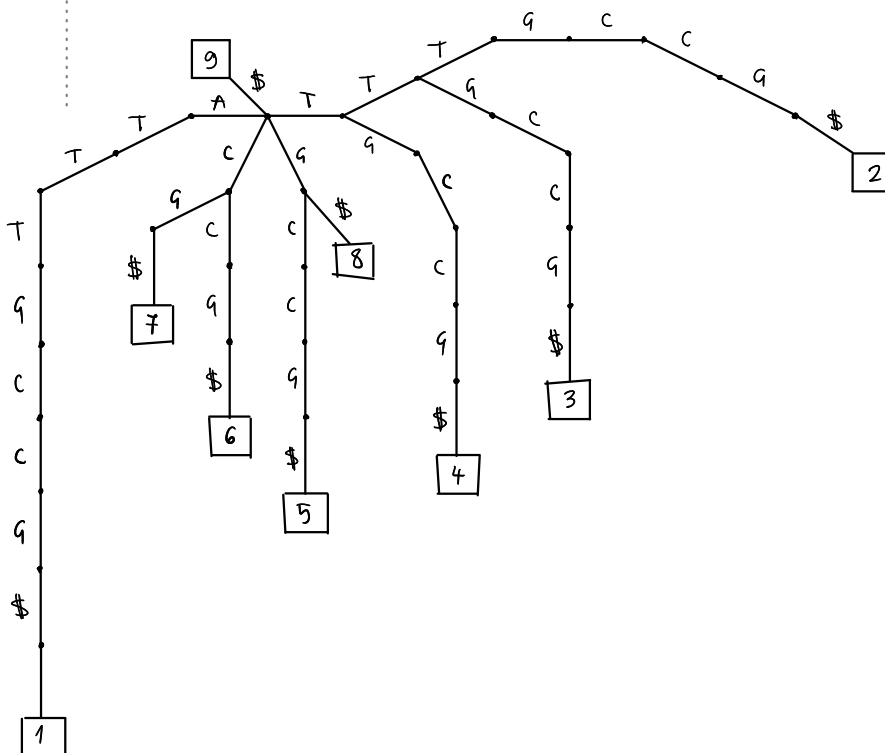
2. Izgradiš drevo v vrstnem redu (abecednem):

1	A	T	T	T	G	C	C	G	\$
2	T	T	T	G	C	C	G	\$	
3	T	T	G	C	C	G	\$		
4	T	G	C	C	G	\$			
5	G	C	C	G	\$				
6	C	C	G	\$					
7	C	G	\$						
8	G	\$							
9	\$								

- a) Spustimo se po drevesu do listov; če je razvejitev, vzamemo za rezultat vse možne liste.

Za TT je to indeks 2 in 3.

- b) Za CGT je to indeks ... ni ga.



3. Zgradite priponsko polje za tekst iz naloge 2 in odgovorite na enaki vprašanji kot v nalogi 2.

Indekse listov drevesa zapišemo v polje po vrstnem redu:
Za iskanje vzorca uporabimo bisekcijo.

9	1	7	6	8	5	4	3	2
1	2	3	4	5	6	7	8	9

- a) "TT" 1. Srednji indeks je 5, list/pripona 8 = "G\$". Primerjamo vzorec: $G\$ < TT \Rightarrow$ glej desno.
2. Srednji indeks je 4 (vzamemo levega): $TG... < TT \Rightarrow$ desno.

9	1	7	6	8	5	4	3	2
5	6	7	8	9				

G je "manjši" od T

3.

9	1	7	6	8	5	4	3	2
							8	9
9	1	7	6	8	5	4	3	2
								9

Srednji indeks je 3: $TTG... = TT \Rightarrow$ ujemanje, ampak glej naprej; desno.

Indeks 2: $TTT... = TT \Rightarrow$ ujemanje! Zmanjkalo polja.

b) Analogno.

4. Dano imamo besedilo $t = \text{ATTTGCCG}$.

- (a) Zakodirajte besedilo z uporabo Huffmanovega algoritma.
 (b) Zakodirajte besedilo z uporabo Burrows-Wheeler kompresije.
 (c) Iz dobljene kode pridobite originalen text t .

b) ① BWT ② MTF ③ HE

①

1	A	T	T	T	G	C	C	G	\$
2	T	T	T	G	C	C	G	\$	T
3	T	T	G	C	C	G	\$	A	T
4	T	G	C	C	G	\$	A	T	T
5	G	C	C	G	\$	A	T	T	T
6	C	C	G	\$	A	T	T	T	G
7	C	G	\$	A	T	T	T	G	C
8	G	\$	A	T	T	T	G	C	C
9	\$	A	T	T	T	G	C	C	G

Pripono razvrstimo pa abecednem vrstnem redu. Zadnji stolpec je rezultat naše transformacije:

9	\$...	G		
1	A	...	\$		
6	C	C	...	G	
7	C	G	...	C	
8	G	\$...	C	
5	G	C	...	T	
4	T	G	...	T	
3	T	T	G	...	T
2	T	T	T	...	T

Po vrsti BWT(t) vzemamo znake in jih dajemo na vrh, ostale premaknemo dol, da zapolnimo luknjo.

②

1	\$	G	\$	G	C	C	T	T	T
2	A	\$	G	\$	G	:	C	:	:
3	C	A	A	A	\$		G		
4	G	C	C	C	A		\$		
5	T	T	T	T	A				
		4	2	2	4	1	5	1	1

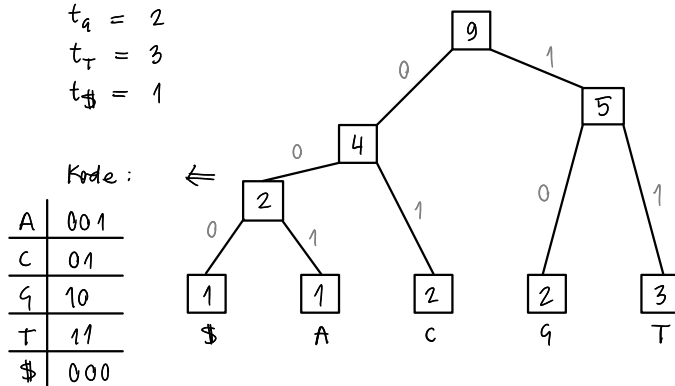
Vrstica iz katere smo vzeli znak.

$\Rightarrow \text{MTF}(\text{BWT}(t)) = 42241511$

zabili smo se enega znaka \Rightarrow krajše kodiranje!

a) $t_A = 1$
 $t_C = 2$
 $t_G = 2$
 $t_T = 3$
 $t_\$ = 1$

12 frekvenc izgradimo drevo.

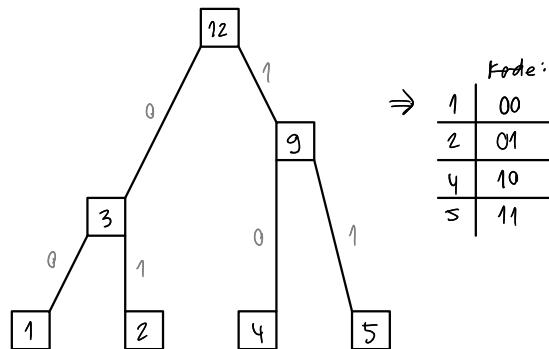


Št. bitov po kodiranju: $t_A \times |\text{HE}(A)| + \dots$
 $= t_A \times 3 + t_C \times 2 + t_G \times 2 + t_T \times 2 + t_\$ \times 3 = 3 + 4 + 4 + 6 + 3 = 20\text{-bitov}$

$\Rightarrow \text{BWT}(t) = G\$GCC TTT$

Zdaj je v lepši obliki za stiskanje.

③ $t_1 = 3$
 $t_2 = 2$
 $t_4 = 2$
 $t_5 = 1$



Št. bitov po kodiranju: $(3+2+2+1) \times 2 = 8 \times 2 = 16\text{-bitov}$

4 manj kot v (a)!

c) encoded(t) = 1 0 0 1 0 1 1 0 0 0 1 1 0 0 0 0
 4 2 2 4 1 5 1 1

Uporabimo MTF tabelo

Ali pa pač preberemo prvo vrstico.

1	\$	G	\$	G	C	C	T	T	T
2	A	\$	G	\$	G	G	C	C	C
3	C	A	A	A	\$	\$	G	G	G
4	G	C	C	C	A	A	\$	\$	\$
5	T	T	T	T	T	A	A	A	
		4	2	2	4	1	5	1	1

\Rightarrow dobimo $G\$GCC TTT$.

\Rightarrow Iz tabele po vrsti glede na številke vzemamo prve znake besede:

$G \$ T T T T G C C G$
 1 2 3 4 5 6 7 8 9
 $A T T T G C C G \$$
 Originalno besedilo.

9	\$...	G		
1	A	...	\$		
6	C	C	...	G	
7	C	G	...	C	
8	G	\$...	C	
5	G	C	...	T	
4	T	G	...	T	
3	T	T	G	...	T
2	T	T	T	...	T

Pove kaj izberemo v prejšnjem stolpcu