

8. Burrows-Wheeler transzformáció

A BURROWS-WHEELER transzformáció  $n$  hosszú szöveget képez le  $(n + 1)$  hosszú szövegre, visszafordíthatóan. A  $T[1..n]$  szöveg transzformációja három lépésben végezhető el: (1)  $T$  végére biggyesztünk egy speciális # karaktert, ami minden más szövegbeli szimbólumtól különbözik, és azok előtt van az ábécében; (2) tekintjük az  $M$  mátrixot, aminek sorai a  $T\#$  ciklikus eltoljtait tartalmazzák, lexikografikusan növekvő sorba rendezve; (3) az  $M$  mátrix utolsó oszlopa adja  $T$  Burrows-Wheeler transzformáltját (bwt). A Burrows-Wheeler transzformáció ekvivalens a suffix array (SA) reprezentációval, ami a  $T\#$  szuffixainak lexikografikus sorrendjét adja meg, az  $1..n + 1$  startpozíciók permutációjaként.

1	l	l	a	l	a	l	a	l	o	12	#	l	a	l	a	l	a	l	o
2	a	l	a	l	a	l	a	l	l	2	a	l	a	l	a	l	a	l	l
3	l	a	l	a	l	a	l	l	a	5	a	l	a	l	a	l	a	l	i
4	i	a	l	a	l	a	l	l	a	9	a	l	a	l	a	l	a	l	l
5	a	l	a	l	a	l	l	a	l	4	i	a	l	a	l	a	l	a	l
6	l	a	l	a	l	a	l	l	a	7	i	a	l	a	l	a	l	a	l
7	i	a	l	a	l	a	l	l	a	1	l	a	l	a	l	a	l	a	#
8	l	a	l	a	l	a	l	l	a	8	l	a	l	a	l	a	l	a	i
9	a	l	a	l	a	l	l	a	l	3	l	a	l	a	l	a	l	a	a
10	l	a	l	a	l	a	l	l	a	6	l	a	l	a	l	a	l	a	a
11	o	#	l	a	l	a	l	a	l	10	l	a	l	a	l	a	l	a	a
12	#	l	a	l	a	l	a	l	a	11	o	#	l	a	l	a	l	a	l

⇒

bwt(lalialilalo#)	=	o	l	i	l	l	l	#	i	a	a	a	l
SA(lalialilalo#)	=	12, 2, 5, 9, 4, 7, 1, 8, 3, 6, 10, 11											

A gyakorlatban az  $M$  mátrixot nem tároljuk sehhol (hiszen az  $O(n^2)$  memóriát igényelne), de egy  $T[1..n]$  szöveg bwt( $T\#$ ) Burrows-Wheeler transzformáltját gyorsan és memória-takarékosan ki lehet számítani<sup>1</sup>.

<sup>1</sup> A bwa illesztő program [Li & Durbin] például létező suffix array implementációt használ, ami SA( $T\#$ )-t  $\Theta(n)$  időben és  $\Theta(n)$  memóriával számolja.

**Tétel.** Minden szövegre bwt( $T\#$ ) egyértelműen meghatározza az  $M$  konceptuális mátrixot, és így az eredeti  $T$  szöveget is.

A BWT egyik előnye, hogy jól tömöríthető akármilyen ábécére (a bzip is erre épül), mivel repetitív szöveg transzformáltjában gyakran hosszan ismétlődik ugyanaz a karakter. DNS szekvencia BWT-jét azonban egyszerűen, direktben lehet tárolni 2 biten pozíciónként.

LF mapping

A BWT-vel való kereséshez azt a tulajdonságát használjuk ki, hogy az azonos karakterek sorrendje az utolsó oszlopban megegyezik a nekik megfelelő karakterek sorrendjével az első oszlopban. Azaz, ha  $LF(i)$  azt jelöli, hogy az  $M$  mátrix  $i$ -edik sorának utolsó karaktere az első oszlopban melyik sorban van, akkor minden  $M[i, n + 1] = M[i', n + 1]$  esetén  $LF(i) > LF(i')$  akkor és csak akkor ha  $i > i'$ . Legyen  $B = \text{bwt}(T\#)$ . Az LF-mapping számolásához kétféle adatot tárolunk:

#l  
a  
a  
a  
a  
i  
l  
l  
l  
l  
l  
l  
l  
l  
l  
l  
l  
o

★  $C[a]$  az  $a$  karakternél kisebb szimbólumok  $B$ -beli számát tartalmazza

★  $\text{Occ}(i, a)$  az  $a$  karakter előfordulását számolja a  $B[1..i]$  prefixben

Az LF-mapping számításához ekkor  $LF(i) = C[B[i]] + \text{Occ}(i, B[i])$ . Ezt a BWT invertálásához is lehet használni: legyen  $s \leftarrow 1$  és  $T[n] \leftarrow B[1]$ , aztán  $i \leftarrow n - 1, n - 2, \dots, 1$  ciklusban legyen  $s \leftarrow LF(s)$  és  $T[i] \leftarrow B[s]$ .

	#	a	i	l	o	
C	0	1	4	6	11	
Occ	a	i	l	o	#	
1	0	0	0	1	0	
2	0	0	1	1	0	
3	0	1	1	1	0	
4	0	1	2	1	0	
5	0	1	3	1	0	
...						
12	3	2	5	1	1	

*Keresés*

Tegyük fel, hogy a  $P[1..p]$  szó előfordulását akarjuk megszámolni a  $T$ -ben. Ehhez az *LF-mapping*-et használjuk: egy  $i \rightarrow p, p-1, \dots, 1$  ciklusban a  $P[i..p]$  szuffixek előfordulását követjük az  $M$  mátrix soraiban.

```

COUNTOCCURRENCE( $P[1..p]$ )
S1  $a \leftarrow P[p]; i \leftarrow p$ 
S2  $s \leftarrow C[a] + 1; e \leftarrow C[a + 1]$  // $P[p]$ -t megtaláljuk az  $s..e$  sorokban
S3 while  $s \leq e$  &&  $i > 1$  do
S4    $i \leftarrow i - 1; a \leftarrow P[i]$ 
S5    $s \leftarrow C[a] + \text{Occ}(s - 1, a) + 1; e \leftarrow C[a] + \text{Occ}(e, a)$ 
      //most  $P[i..p]$ -t találjuk az  $s..e$  sorokban
S6 return  $\max\{e - s + 1, 0\}$ 

```

*BWA*

Tegyük fel, hogy nem megszámolni akarjuk, hanem azt tudni, hogy milyen pozíciókban fordul elő  $P$ . A keresés hasonlóan zajlik le, mint az előfordulások számolása, de most a megtalált sorokhoz tartozó  $SA(s)$ ,  $SA(s + 1)$ ,  $\dots$ ,  $SA(e)$  értékeket kell tudni. Ezt akár a BWT tömörített változatával is meg lehet csinálni nagyon hatékonyan. Az általános ábécére működő kompakt adatstruktúrát FM-indexnek hívják, Paolo Ferragina és Giovanni Manzini fejlesztette ki. DNS-re azonban lehet ennél praktikusabb megoldás is. A bwa program [Li & Durbin] a Burrows-Wheeler transzformálton kívül ( $2n$  bit) az Occ értékeket közvetlenül tárolja a négy nukleotidra, minden 128. pozícióban ( $n/32$  integer változó), és a SA értékeket szintén, minden 32. pozícióban ( $n/32$  integer változó). Köztes  $i$  pozíciókban az SA értéke kiszámolható az LF-mapping segítségével, Occ pedig a nukleotidok direkt számolásával. Az adatstruktúrát tárolni kell mind a forward mind a reverse szárra. Mindösszesen néhány Gigabyte elég az emberi genom teljes indexeléséhez.

A Burrows-Wheeler alapú illesztés nehézsége, hogy míg pontos előfordulást nagyon gyorsan megtalál, a mismatch-ek és a gap-ek kezelése elég nehézkes. Ehhez a bowtie [Langmead ++] és a bwa program is úgy keres, hogy max.  $d$  hibát megenged a keresett szóban, és a hibákkal módosított variánsokra is rákeres.

*Hivatkozások*

M. Burrows and D. J. Wheeler. A block-sorting lossless data compression algorithm. Technical Report SRC 124, Digital Equipment Corporation, 1994. <http://www.hpl.hp.com/techreports/Compaq-DEC/SRC-RR-124.pdf>.

Paolo Ferragina and Giovanni Manzini. Opportunistic data structures with applications. In *Foundations of Computer Science (FOCS)*, pages 390–398, 2000. <http://people.unipmn.it/manzini/papers/focs00draft.pdf>.

Ben Langmead, Cole Trapnell, Mihai Pop, and Steven L. Salzberg. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biology*, 10:R25, 2009. DOI: 10.1186/gb-2009-10-3-r25.

Heng Li and Richard Durbin. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, 25:1754–1760, 2009. DOI: 10.1093/bioinformatics/btp324.