

# Komprese dat pomocí transformace Burrows–Wheeler

Petr Dvořáček – xdvora0n@stud.fit.vutbr.cz

6. května 2014

Komprese pomocí Burrows–Wheelerovy transformace se řadí mezi blokové komprese. Tento algoritmus nachází využití i v bioinformatice pro sestavování genomu organismů z krátkých fragmentů (řetězců) DNA.

Tato práce se však zabývá prvním problémem a tou je komprese dat. Samotná komprese je rozdělena do několika kroků. Nejdříve se vstupní text ztransformuje pomocí Burrows–Wheelerovy transformace (dále jen BWT), jak je uvedeno v sekci 1. Výstup BWT zpracujeme pomocí funkce Move to Front (dále jen MFT). O tomto algoritmu pojednává sekce 2. Move to Front zvyšuje efektivnost entropického kódování. Data získána MTF jsou pak vstupem do algoritmu Run Length Encoding, který provádí zkrácení posloupnosti stejných bytů. Vizte sekci 3. Výstup RLE je nakonec zpracován pomocí entropického kódování. V tomto projektu jsem použil statické Huffmanovo kódování. O tom pojednávám v sekci 4. Následují experimenty a závěr.

## 1 Burrows–Wheelerova transformace

Burrows wheelerova transformace vyrábí permutaci znaků z řetězce  $x$  o velikosti  $n$ . Snaží se znaky seřadit tak, abychom získali velkou posloupnost stejných znaků za minimální cenu (z pohledu paměti).

Nejdříve řetězec  $x$ , kde například  $x = \text{SWISS\_MISS}$ , se uspořádá do matice  $M = n \times n$ . Na jednotlivých řádcích aplikujeme operátor rotace původního řetězce. Tyto rotované řádky se pak seřadí lexiograficky. Výstupem je pak poslední sloupec s indexem řádku  $i$  na původní řetězec  $x$ . Na příkladu je to znázorněno divně modrou barvou, kde  $f(x) = \text{SWM\_SISSS}$  a  $i = 8$ .

Původní text získáme následujícím způsobem. Jednotlivá písmena se seřadí stabilním řadícím algoritmem. První písmeno původního textu se nachází v seřazeném poli na uložené pozici  $i$ . Nový index je pak hodnota indexu nalezeného písmena v neseřazeném poli. Ten pak určuje pozici druhého písmena v seřazeném poli. Analogicky získáme původní řetězec, tento postup je znázorněn níže:

0	SWISS_MISS	5	_MISSSWISS	0	S	0→3	-	0→3	-	0→3	-
1	WISS_MISSS	2	ISS_MISSSSW	1	W	1→5	I	1→5	I	1→5	I
2	ISS_MISSSSW	7	ISSSWISS_M	2	M	2→6	I	2→6	I	2→6	I
3	SS_MISSSSWI	6	MISSSWISS_	3	-	3→2	M	3→2	M	3→2	M
4	S_MISSSSWIS	4	S_MISSSSWIS	4	S	4→0	S	4→0	S	4→0	S
5	_MISSSWISS	3	SS_MISSSSWI	5	I	5→4	S	5→4	S	5→4	S
6	MISSSWISS_	8	SSSWISS_MI	6	I	6→7	S	6→7	S	6→7	S
7	ISSSWISS_M	9	SSWISS_MIS	7	S	7→8	S	7→8	S	7→8	S
8	SSSWISS_MI	0	SWISS_MISS	8	S	8→9	S	8→9	S	8→9	S
9	SSWISS_MIS	1	WISS_MISSS	9	S	9→1	W	9→1	W	9→1	W

  

0→3	-	0→3	-	0→3	-	0→3	-	0→3	-	0→3	-
1→5	I	1→5	I	1→5	I	1→5	I	1→5	I	1→5	I
2→6	I	2→6	I	2→6	I	2→6	I	2→6	I	2→6	I
3→2	M	3→2	M	3→2	M	3→2	M	3→2	M	3→2	M
4→0	S	4→0	S	4→0	S	4→0	S	4→0	S	4→0	S
5→4	S	5→4	S	5→4	S	5→4	S	5→4	S	5→4	S
6→7	S	6→7	S	6→7	S	6→7	S	6→7	S	6→7	S
7→8	S	7→8	S	7→8	S	7→8	S	7→8	S	7→8	S
8→9	S	8→9	S	8→9	S	8→9	S	8→9	S	8→9	S
9→1	W	9→1	W	9→1	W	9→1	W	9→1	W	9→1	W

## 2 Move to Front

Myšlenkou této části je, že se vytvoří seznam naposledy použitých symbolů. Vygeneruje se takovým způsobem, že bude obsahovat 256 položek, s hodnotami  $(0, 1, 2, 3, \dots, 254, 255)$ . Seznam se aktualizuje tak, že se na jeho počátek dá element se vstupní hodnotou symbolu. Výstupem je pak index, kde byla hodnota nalezena. Inverzní operace je podobná s tím rozdílem, že vstupem je index do slovníku a výstupem je hodnota symbolu.

Příklad: Nechť máme seznam o 26 položkách s hodnotami  $(A \dots Z)$  a vstupní řetězec  $x = \text{SWMSIISSS}$ .

0. iterace	$(A, B, C, D, \dots, Z)$	S	W	M	S	I	I	S	S	S
1. iterace	$(S, A, B, C, \dots, Z)$	18	W	M	S	I	I	S	S	S
2. iterace	$(W, S, A, B, \dots, Z)$	18	22	M	S	I	I	S	S	S
3. iterace	$(M, W, S, A, \dots, Z)$	18	22	12	S	I	I	S	S	S
4. iterace	$(S, M, W, A, \dots, Z)$	18	22	12	2	I	I	S	S	S
...	...									
9. iterace	$(S, I, M, W, \dots, Z)$	18	22	12	2	8	0	1	0	0

## 3 Run Length Encoding

Tato část algoritmu se snaží komprimovat posloupnosti shodných znaků. Tato posloupnost je zakódována do formátu *délka, znak*. Můžeme dosáhnout lepší komprese dat, zmenšujeme-li posloupnosti délky větší než dva. Pak použijeme formát *speciální znak, délka, znak* s důmyslným ošetřením výskytu speciálního znaku, např. jeho zdvojení.

Mějme například řetězec  $x = \# \text{BAAANAAAAANAAAAA}$ , kde speciální znak je  $\#$ . Pak při použití RLE získáme posloupnost:  $rle(x) = \# \# \text{B} \# 3 \text{AN} \# 4 \text{AN} \# 9 \text{A}$ .

## 4 Statické Huffmanovo kódování

V tomto projektu bylo použito statické Huffmanovo kódování. Huffmanův strom byl sestaven podle tabulky uvedné níže. Předpokládám, že nejvíce zastoupené symboly nesou malé hodnoty.

Listů	Hloubka	Prvky
2	2	0, 1
4	4	2, 3, 4 + 10 ta měla vyšší pravděpodobnost než 5 a 4.
8	6	5, 6, 7, 8, 9, 11, 12 + speciální znak 255
16	8	13, 14 atd. + maximální délka RLE 254
32	10	28, 29 ..
64	12	60, 61 ..
127	13	124, 125 ..
4	15	Prvky 251, 252, 253 + speciální znak konce bloku.

## 5 Experimentální vyhodnocení a závěr

Projekt byl naimplementován v jazyce C a byl otestován na školním serveru Merlin (škoda, že se nejmenuje Gandalf). Velikost bloku byla nastavena na 200000. Pro změnu stačí přeložit příkazem: `make BLOCK_SIZE=42`. Podařilo se komprimovat 751kB referenční soubor na 301kB.

Další práci do budoucna bych viděl v optimalizaci programu. Některé části (např. MTF) by se daly provést efektivnějším způsobem (např. pomocí `memcpy()`). Rovněž není možnost využití megabytových bloků, neboť by byla pro ně nutná alokace pomocí funkce `malloc()`.