

bwt, v. 0.2: Compute the Burrows-Wheeler Transform

Bernhard Haubold

Max-Planck-Institute for Evolutionary Biology, Plön, Germany

March 9, 2012

1 Introduction

The Burrows-Wheeler Transform (BWT) [3] is an integral part of many compression algorithms, including the widely used program `bzip2` for compacting files. In Bioinformatics the BWT underlies a number of highly memory-efficient tools, including `bwa` [7] and `bowtie` [6].

My program `bwt` demonstrates the encoding and the decoding phase of the BWT. I had two aims when writing it: (i) to demonstrate the BWT when teaching Bioinformatics, and (ii) to explore the properties of the BWT when applied to real-world data sets ranging from Shakespeare's *Hamlet* to bacterial genomes and proteomes.

In the next Section I explain how to get started with `bwt`. This is followed by a tutorial-style exposition of the central ideas behind the BWT as reflected in the features of `bwt`. My approach is heavily indebted to [1], which ought to be consulted for more details on the BWT and its relationship to other string-centered data structures, especially suffix trees and suffix arrays.

2 Getting Started

`bwt` was written in C on a computer running Linux and should work on any standard UNIX system. However, please contact me at `haubold@evolbio.mpg.de` if you have any problems with the program.

- Unpack the program

```
tar -xvzf bwt_XXX.tgz
```

where XXX indicates the version.

- Change into the newly created directory

```
cd Bwt_XXX
```

and list its contents

```
ls
```

- Generate `bwt`

```
make
```

- List its options

```
./bwt -h
```

- `bwt` takes FASTA-formatted input, for example

```
./bwt Data/mississippi.fasta
```


A		B		C	
\mathcal{F}	\mathcal{L}	\mathcal{F}	\mathcal{L}	\mathcal{F}	\mathcal{L}
\$	i	\$ ₁	i ₁	\$ ₁	i ₁
i	p	i ₁	p ₁	i ₁	p ₁
i	s	i ₂	s ₁	i ₂	s ₁
i	s	i ₃	s ₂	i ₃	s ₂
i	m	i ₄	m ₁	i ₄	m ₁
m	\$	m ₁	\$ ₁	m ₁	\$ ₁
p	p	p ₁	p ₂	p ₁	p ₂
p	i	p ₂	i ₂	p ₂	i ₂
s	s	s ₁	s ₃	s ₁	s ₃
s	s	s ₂	s ₄	s ₂	s ₄
s	i	s ₃	i ₃	s ₃	i ₃
s	i	s ₄	i ₄	s ₄	i ₄

Figure 1: Decoding the Burrows-Wheeler Transform shown in column \mathcal{L} . **A**: Pairing the transform with its sorted version. **B**: Counting the occurrences of each character. **C**: Lats-first mapping starting from $\$$ ₁ in \mathcal{L} to recover the input.

But what feature of *Hamlet* induces this run of n? To find out, we print the context of the same segment:

```
./bwt -c -p 91554-92113 Data/hamlet.fasta | less
```

In other words, because *g* is frequently preceded by *n*, the transform groups *ns* by sorting on *g*.

3.2 Decoding

The method for retrieving the original string from its BWT is surprisingly simple: Consider again the original string *mississippi\$*, which we have transformed to *ipssm\$piissii*. Remember that this is the last column, \mathcal{L} , of the sorted matrix of string rotations. We write \mathcal{L} next to the first column of this matrix, \mathcal{F} (Figure 1A). Recall that \mathcal{F} is just the alphabetically ordered input string. The trick is now to count the occurrences of each character in \mathcal{F} and \mathcal{L} such that c_i is the *i*-th occurrence of *c* (Figure 1B). To recover the input string, we start at the position of the sentinel character, $\$$ ₁ in \mathcal{L} ; the corresponding character in \mathcal{F} is the first character, *m*₁, of the input string. Next, search for *m*₁ in \mathcal{L} and look up its partner in \mathcal{F} , and so on (Figure 1C). The end result of this single table traversal is *m₁i₄s₄s₂i₃s₃s₁i₂p₂p₁i₁\$₁*.

In practice it is not necessary to explicitly construct \mathcal{F} . Let *T* be the text we consider; then three auxiliary arrays suffice to reconstruct \mathcal{F} from \mathcal{L} :

- $C[i]$: the number of times the character $T[i]$ has occurred before position *i*;
- $K[c]$: the total count of character *c*;
- $M[c]$: the position at which character *c* first appears in \mathcal{F} .

Decoding in demo mode returns all this information (and a bit more):

```
./bwt -d -D Data/mississippiEncoded.fasta
      K:      M:
P: 123456789012   $ 1      $ 1      P: 123456789012
S: ipssm$piissii  i 4      i 2      S: mississippi$
C: 000100112323  m 1      m 6
                   p 2      p 7
                   s 4      s 9
```

The details of how these arrays are used to reconstruct the input string from \mathcal{L} are explained in [1, chp. 2]. More important than these details is the fact that decoding only takes time proportional to the length of the text. This

contrasts with encoding, which is based on sorting n strings of length n . Sorting usually takes time proportional to $n \log(n)$.

To convince ourselves that decoding really reverses the work of encoding, try the following:

```
./bwt Data/hamlet.fasta | bwt -d > tmp.fasta
diff tmp.fasta Data/hamlet.fasta
1c1
< >hamlet - bwt - decoded
---
> >hamlet
```

Only the header lines differ, but the text is unchanged after one round of encoding and decoding.

Instead of *Hamlet*, we can also subject the genome sequence of the standard laboratory strain of *Escherichia coli* to BWT:

```
./bwt Data/k12genome.fasta | less
```

A cursory glance at the transform reveals no obvious long runs of identical nucleotides.—The sequence looks like a “normal” genome. We can quantify this impression of low repetitiveness by computing the index of repetitiveness, I_r , using the program *ir* [4], which is freely available from my home page under the link “Software”.

```
./bwt Data/k12genome.fasta | ir
# Len I_r
4639676 0.0010
```

In fact, the I_r for the transform is much lower than for the original genome:

```
ir < Data/k12genome.fasta
# Len I_r
4639675 0.0587
```

In contrast, the original *Hamlet* has a lower I_r than its transform:

```
ir < Data/hamlet.fasta
# Len I_r
176686 -0.0174
```

```
./bwt Data/hamlet.fasta | ir
# Len I_r
176687 0.0430
```

This would indicate that in natural languages a given suffix determines much more strongly the preceding character than in genome sequences.

What about the *E. coli* proteome? We begin by converting it into a single string:

```
echo '>k12proteome' > proteome.fasta
grep -v '^>' Data/k12proteome.fasta >> proteome.fasta
```

The repetitiveness of this string is substantial:

```
ir < proteome.fasta
# Len I_r
1297495 0.3479
```

However, as with the genome, the BWT *reduces* repetitiveness:

```
./bwt proteome.fasta | ir
# Len I_r
1297496 -0.0488
```

We can compare this to the I_r of the randomized proteome using `randomizeSeq`, which is available from my homepage under “Software→bioBox”:

```
randomizeSeq proteome.fasta | ir
# Len I_r
1297495 -0.0496
```

The I_r of the randomized proteome is almost identical to that of the transform. This illustrates that the dependence between consecutive amino acids is so weak that the BWT effectively shuffles a proteome. The independence of adjacent amino acids was already observed in insulin, the very first protein sequenced. At the time, the free association between amino acids led to the rejection of Crick’s commaless genetic code [5]. The essential incompressibility of protein sequences has repeatedly been investigated since then [2]. So it is less the compression aspect that motivates application of the BWT in Bioinformatics than the fact that the BWT is a space-efficient replacement of enhanced suffix arrays, which in turn are space efficient representations of suffix trees [1].

4 Listings

4.1 The Driver Program, `bwt.c`

```
1  /***** bwt.c *****/
   * Description: Program for exploring the Burrows-
   *   Wheeler Transform.
   * Reference: D. Adjero, T. Bell, and A.
   *   Mukherjee (2008). The Burrows-Wheeler Trans-
6  *   form; Data Compression, Suffix Arrays, and
   *   Pattern Matching. Springer.
   * Author: Bernhard Haubold, haubold@evolbio.mpg.de
   * Date: Fri Mar 2 08:18:05 2012
   *****/
11 #include <stdio.h>
   #include <stdlib.h>
   #include <unistd.h>
   #include <fcntl.h>
   #include "interface.h"
16 #include "eprintf.h"
   #include "sequenceData.h"

   void scanFile(int fd, Args *args);
   void decode(Args *args, Sequence *seq);
21 void encode(Args *args, Sequence *seq);

   int main(int argc, char *argv[]){
       int i;
       char *version;
26   Args *args;
       int fd;

       version = "0.2";
       setprogname2("bwt");
31   args = getArgs(argc, argv);
       if(args->h || args->e)
           printUsage(version);
       if(args->v)
           printSplash(version);
```

```

36  if(args->numInputFiles == 0){
        fd = 0;
        scanFile(fd, args);
    }else{
        for(i=0;i<args->numInputFiles;i++){
41      fd = open(args->inputFiles[i],0);
        scanFile(fd, args);
        close(fd);
        }
    }
46  free(args);
    free(progname());
    return 0;
}

```

```

51 void scanFile(int fd, Args *args){
    Sequence *seq;

    seq = readFasta(fd,args->s);
    if(args->d)
56  decode(args, seq);
    else
        encode(args, seq);
    freeSequence(seq);
}

```

4.2 Encoding, encoding.c

```

/***** encode.c *****
* Description: Encoding step of Burrows-Wheeler
* Transform.
* Reference: Donald Adjero, Timorhy Bell, and
5 * Amar Mukherjee (2008). The Burrows-Wheeler
* Transform; Data Compression, Suffix Arrays,
* and Pattern Matching. Springer.
* Author: Bernhard Haubold, haubold@evolbio.mpg.de
* Date: Mon Mar 5 11:58:31 2012
10 *****/
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
#include "interface.h"
15 #include "eprintf.h"
#include "sequenceData.h"

char *globalText;
int globalTextLen;

20 int ringStrCmp(const void *v1, const void *v2);
void printRotations(char *text, int *indexArr, int n);
void printTransform(Args *args, char *text, char *header, int *indexArr,
    int n);
void printContext(Args *args, char *text, char *header, int *indexArr, int
    n);
25 void printRuns(Args *args, char *text, char *header, int *indexArr, int n);

```

```

void encode(Args *args, Sequence *seq){
    int *indexArr;
    int i;
30
    indexArr = (int *)emalloc(seq->len*sizeof(int));
    for(i=0;i<seq->len;i++)
        indexArr[i] = i;
    globalTextLen = seq->len;
35
    globalText = seq->seq;
    if(args->D){
        printf("Text:\n");
        printf("%s\n\n",seq->seq);
        printf("Rotations:\n");
40
        printRotations(globalText, indexArr, globalTextLen);
    }
    qsort(indexArr,globalTextLen,sizeof(int),ringStrCmp);
    if(args->D){
        printf("\nSorted Rotations:\n");
45
        printRotations(globalText, indexArr, globalTextLen);
        printf("\nTransform:\n");
        printTransform(args, globalText, NULL, indexArr, globalTextLen);
    }
    if(!args->D){
50
        if(!args->c){
            if(args->r)
                printRuns(args, globalText, seq->headers[0], indexArr,
                    globalTextLen);
            else
                printTransform(args, globalText, seq->headers[0], indexArr,
                    globalTextLen);
55
        }else
            printContext(args, globalText, seq->headers[0], indexArr,
                globalTextLen);
    }

    free(indexArr);
60 }

void printContext(Args *args, char *text, char *header, int *indexArr, int
    n){
    int i, j, start, end;
65
    if(args->p){
        start = args->p[0] > -1 ? args->p[0] : 0;
        end = args->p[1] + 2 < n ? args->p[1] + 2 : n;
    }else{
        start = 0;
70
        end = n+1;
    }

    for(i=start;i<end-1;i++){
        printf("%d:␣",i+1-start);
75
        for(j=indexArr[i];j<indexArr[i]+args->l;j++)

```

```

        printf("%c",text[j%n]);
        printf("\n");
    }
80 }

void printRuns(Args *args, char *text, char *header, int *indexArr, int n){
    int i, j, o, start, end;
    int rowLen = 70;
85 char c1, c2;
    char *count;

    count = (char *)emalloc(sizeof(char)*256);
    for(i=0;i<256;i++)
90     count[i] = 0;

    if(header)
        printf("%s_\_bwt\n",header);
    o = n-1;
95 j = 0;
    if(args->p){
        start = args->p[0] > -1 ? args->p[0] : 0;
        end = args->p[1] + 2 < n ? args->p[1] + 2 : n;
    }else{
100     start = 0;
        end = n+1;
    }
    c1 = text[(indexArr[start]+o)%n];
    count[(int)c1] = 1;
105 for(i=start+1;i<end-1;i++){
        c2 = text[(indexArr[i]+o)%n];
        if(c2 != c1){
            if(count[(int)c1] > 1)
                printf("%c%d",c1,count[(int)c1]);
110             else
                printf("%c",c1);
            j += (log(count[(int)c1])/log(10)+1) + 1;
            count[(int)c1] = 0;
            if(j>=rowLen){
115                 j = 0;
                printf("\n");
            }
        }
        count[(int)c2]++;
120     c1 = c2;
    }
    if(count[(int)c1] > 1)
        printf("%c%d",c1,count[(int)c1]);
    else
125     printf("%c",c1);
    printf("\n");
    free(count);
}

```



```

130 void printTransform(Args *args, char *text, char *header, int *indexArr,
    int n){
    int i, j, o, start, end;
    int rowLen = 70;

135    if(header)
        printf("%s_\bwt\n",header);
    o = n-1;
    j = 0;
    if(args->p){
140        start = args->p[0] > -1 ? args->p[0] : 0;
        end = args->p[1] + 2 < n ? args->p[1] + 2 : n;
    }else{
        start = 0;
        end = n+1;
145    }

    for(i=start;i<end-1;i++){
        j++;
        printf("%c",text[(indexArr[i]+o)%n]);
150        if(j==rowLen){
            j = 0;
            printf("\n");
        }
    }
155    if(j)
        printf("\n");
}

void printRotations(char *text, int *indexArr, int n){
160    int i, j;

    for(i=0;i<n;i++){
        for(j=0;j<n;j++){
            printf("%c",text[(indexArr[i]+j)%n]);
165            printf("\n");
        }
    }

170 int ringStrCmp(const void *v1, const void *v2){
    int a, b, i;

    a = *(int *)v1;
    b = *(int *)v2;
175    i = 0;

    while(globalText[a % globalTextLen] == globalText[b % globalTextLen] && i
        < globalTextLen){
        a++;
        b++;
180        i++;
    }
}

```

```

    if(globalText[a] == globalText[b])
        return 0;
185 else if(globalText[a] < globalText[b])
        return -1;
    else
        return 1;
}

```

4.3 Decoding, decoding.c

```

/***** decode.c *****/
* Description: BWT-decode.
* Reference: Donald Adjero, Timorhy Bell, and
4 *   Amar Mukherjee (2008). The Burrows-Wheeler
*   Transform; Data Compression, Suffix Arrays,
*   and Pattern Matching. Springer, p. 26.
* Author: Bernhard Haubold, haubold@evolbio.mpg.de
* Date: Mon Mar 5 12:22:11 2012
9 *****/
#include <stdio.h>
#include <stdlib.h>
#include "eprintf.h"
#include "sequenceData.h"
14 #include "interface.h"

void decode(Args *args, Sequence *seq){
    int i, j, a, n, sum, dictSize, start, end;
    /* let F be the first and L the last column in the sorted array of
       rotations */
19 int *K; /* K[i]: count of character seq[i] */
int *C; /* C[i]: occurrences of character seq[i] before position i in L
*/
int *M; /* M[i]: F[M[i]]: first time character seq[i] occurs in F */
char *Q; /* output string */

24 dictSize = 256;
n = seq->len - 1;
K = (int *)emalloc(sizeof(int)*dictSize);
M = (int *)emalloc(sizeof(int)*dictSize);
C = (int *)emalloc(sizeof(int)*n);
29 Q = (char *)emalloc(sizeof(char)*n);
/* initialize K */
for(i=0;i<dictSize;i++)
    K[i] = 0;
/* count characters in K */
34 /* record number of previous appearances of character in C */
for(i=0;i<n;i++){
    C[i] = K[(int)seq->seq[i]];
    K[(int)seq->seq[i]]++;
}
39 /* first occurrence of character in F */
sum = 0;
for(i=0;i<dictSize;i++){
    M[i] = sum;

```

```

    sum += K[i];
44 }
/* look for starting character */
for(i=0;i<n;i++)
    if(seq->seq[i] == BORDER){
        a = i;
49     break;
    }
a = i;
for(j=n-1;j>-1;j--){
    Q[j] = seq->seq[i];
54     i = C[i] + M[(int)seq->seq[i]];
}
sum = 0;
if(args->D){
    printf("P:␣");
59     for(i=0;i<n;i++)
        printf("%d", (i+1)%10);
    printf("\nS:␣");
    for(i=0;i<n;i++)
        printf("%c", seq->seq[i]);
64     printf("\nC:␣");
    for(i=0;i<n;i++)
        printf("%d", C[i]);
    printf("\nK:\n");
    for(i=0;i<dictSize;i++)
69     if(K[i]>0)
        printf("%c\t%d\n", i, K[i]);
    printf("M:\n");
    for(i=0;i<dictSize;i++)
    if(K[i]>0)
74     printf("%c\t%d\n", i, M[i]+1);
    printf("P:␣");
    for(i=0;i<n;i++)
        printf("%d", (i+1)%10);
    printf("\nS:␣");
79 }else{
    printf("%s␣-␣decoded\n", seq->headers[0]);
}
if(args->p){
    start = args->p[0] > -1 ? args->p[0] : 0;
84     end = args->p[1] + 2 < n ? args->p[1] + 2 : n;
} else{
    start = 0;
    end = n;
}
89 sum = 0;
for(i=start;i<end-1;i++){
    printf("%c", Q[i]);
    sum++;
    if(sum == 70){
94     printf("\n");
        sum = 0;
    }
}

```

```

    }
    99  if(args->D){
        sum++;
        printf("%c",Q[i]);
    }
    if(sum)
        printf("\n");
104
    free(C);
    free(K);
    free(M);
    free(Q);
109 }

```

5 Change Log

- Version 0.2 (7th March 2012)
 - First released version.

References

- [1] D. Adjeroh, T. Bell, and A. Mukherjee. *The Burrows-Wheeler Transform:: Data Compression, Suffix Arrays, and Pattern Matching*. Springer, 2008.
- [2] D. Adjeroh and Fei Nan. On compressibility of protein sequences. In *Data Compression Conference, 2006. DCC 2006. Proceedings*, pages 10 pp. –434, march 2006.
- [3] M. Burrows and D. J. Wheeler. A block-sorting lossless data compression algorithm. Technical Report 124, Digital Equipment Corporation, Palo Alto, California, 1994.
- [4] B. Haubold and T. Wiehe. How repetitive are genomes? *BMC Bioinformatics*, 7:541, 2006.
- [5] H. F. Judson. *The Eighth Day of Creation*. Penguin Books, London, 1979/1995.
- [6] B. Langmead, C Trapnell, M. Pop, and S. L. Salzberg. Ultrafast and memory-efficient alignment of short dna sequences to the human genome. *Genome Biology*, 10:R25, 2009.
- [7] H. Li and R. Durbin. Fast and accurate short read alignment with burrows-wheeler transform. *Bioinformatics*, 25:1754–1760, 2009.