

Burrows–Wheeler transform

From Wikipedia, the free encyclopedia

The **Burrows–Wheeler transform** (**BWT**, also called **block-sorting compression**) rearranges a character string into runs of similar characters. This is useful for compression, since it tends to be easy to compress a string that has runs of repeated characters by techniques such as move-to-front transform and run-length encoding. More importantly, the transformation is **reversible**, without needing to store any additional data. The BWT is thus a "free" method of improving the efficiency of text compression algorithms, costing only some extra computation.

Contents

■ 1 Description
■ 2 Example
■ 3 Explanation
■ 4 Optimization
■ 5 Bijective variant
■ 6 Dynamic Burrows–Wheeler transform
■ 7 Sample implementation
■ 8 BWT in bioinformatics
■ 9 References
■ 10 External links

Description

The Burrows–Wheeler transform is an algorithm used in data compression techniques such as bzip2. It was invented by Michael Burrows and David Wheeler in 1994 while Burrows was working at DEC Systems Research Center in Palo Alto, California.^[1] It is based on a previously unpublished transformation discovered by Wheeler in 1983.

When a character string is transformed by the BWT, the transformation permutes the order of the characters. If the original string had several substrings that occurred often, then the transformed string will have several places where a single character is repeated multiple times in a row.

For example:

Input	SIX.MIXED.PIXIES.SIFT.SIXTY.PIXIE.DUST.BOXES
Output	TEXYDST.E.IXIXIXXSMPPFS.B..E.S.EUSFXDIOIIIT ^[2]

The output is easier to compress because it has many repeated characters. In this example the transformed string, there are a total of eight runs of identical characters: XX, II, XX, SS, PP, . . . , II, and III, which together make 17 out of the 44 characters.

Example

The transform is done by sorting all rotations of the text into lexicographic order, by which we mean that the 8 rotations appear in the second column in a different order, in that the 8 rows have been sorted into lexicographical order. We then take as output the last column and the number *k* = 7 of the row that the non rotated row ends up in. For example, the text "BANANA" is transformed into "BNN^AA^A" through these steps (the red ^ character indicates the "EOF" pointer):

Transformation				
Input	All rotations	Sorted into lexical order	Taking last column	Output last column
	<div> <div>^"BANANA"</div> <div>^"ANANA"</div> <div>^"NANA"</div> <div>^"ANANA"</div> <div>^"NANA"</div> <div>^"ANANA"</div> <div>^"NANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>
	<div> <div>^"BANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>
	<div> <div>^"BANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>
	<div> <div>^"BANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>
	<div> <div>^"BANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>
	<div> <div>^"BANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>

The following pseudocode gives a simple (though inefficient) way to calculate the BWT and its inverse. It assumes that the input string *s* contains a special character 'EOF' which is the last character, occurs nowhere else in the text.

```

function BWT (string s)
  create a table, rows are all possible rotations of s
  sort rows alphabetically
  return (last column of the table)
-----
function inverseBWT (string s)
  create empty table
  repeat length(s) times
    // first insert creates first column
    insert s as a column of table before first column of the table
    sort rows of the table alphabetically
    return (row that ends with the 'EOF' character)
-----

```

Explanation

To understand why this creates more-easily-compressible data, consider transforming a long English text frequently containing the word "the". Sorting the rotations of this text will group rotations starting with "the" together, and the last character of that rotation (which is also the character before the "he") will usually be "t", so the result of the transform would contain a number of "t" characters along with the perhaps less-common exceptions (such as if it contains "bHrahe") mixed in. So it can be seen that the success of this transform depends upon one value having a high probability of occurring before a sequence, so that in general it needs fairly long samples (a few kilobytes at least) of appropriate data (such as text).

The remarkable thing about the BWT is not that it generates a more easily encoded output—an ordinary sort would do that—but that it is *reversible*, allowing the original document to be re-generated from the last column data.

The inverse can be understood this way. Take the final table in the BWT algorithm, and erase all but the last column. Given only this information, you can easily reconstruct the first column. The last column tells you all the characters in the text, so just sort these characters alphabetically to get the first column. Then, the first and last columns (of each row) together give you all *pairs* of successive characters in the document, where pairs are taken cyclically so that the last and first character form a pair. Sorting the list of pairs gives the first *and second* columns. Continuing in this manner, you can reconstruct the entire list. Then, the row with the "end of file" character at the end is the original text. Reversing the example above is done like this:

Inverse transformation				
Input				
<div> <div>^"BNN^AA^A"</div> </div>				
Add 1	Sort 1	Add 2	Sort 2	
<div> <div>^"B"</div> <div>^"N"</div> <div>^"N"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> </div>	<div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> </div>	<div> <div>^"B"</div> <div>^"N"</div> <div>^"N"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> </div>	<div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> </div>	
Add 3	Sort 3	Add 4	Sort 4	
<div> <div>^"BAN"</div> <div>^"BAN"</div> <div>^"BAN"</div> <div>^"BAN"</div> <div>^"BAN"</div> <div>^"BAN"</div> <div>^"BAN"</div> <div>^"BAN"</div> </div>	<div> <div>^"ANA"</div> <div>^"ANA"</div> <div>^"ANA"</div> <div>^"ANA"</div> <div>^"ANA"</div> <div>^"ANA"</div> <div>^"ANA"</div> <div>^"ANA"</div> </div>	<div> <div>^"BAN"</div> <div>^"BAN"</div> <div>^"BAN"</div> <div>^"BAN"</div> <div>^"BAN"</div> <div>^"BAN"</div> <div>^"BAN"</div> <div>^"BAN"</div> </div>	<div> <div>^"ANAN"</div> <div>^"ANAN"</div> <div>^"ANAN"</div> <div>^"ANAN"</div> <div>^"ANAN"</div> <div>^"ANAN"</div> <div>^"ANAN"</div> <div>^"ANAN"</div> </div>	
Add 5	Sort 5	Add 6	Sort 6	
<div> <div>^"BANAN"</div> <div>^"BANAN"</div> <div>^"BANAN"</div> <div>^"BANAN"</div> <div>^"BANAN"</div> <div>^"BANAN"</div> <div>^"BANAN"</div> <div>^"BANAN"</div> </div>	<div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> <div>^"ANANA"</div> </div>	<div> <div>^"BANANA"</div> <div>^"BANANA"</div> <div>^"BANANA"</div> <div>^"BANANA"</div> <div>^"BANANA"</div> <div>^"BANANA"</div> <div>^"BANANA"</div> <div>^"BANANA"</div> </div>	<div> <div>^"ANANAN"</div> <div>^"ANANAN"</div> <div>^"ANANAN"</div> <div>^"ANANAN"</div> <div>^"ANANAN"</div> <div>^"ANANAN"</div> <div>^"ANANAN"</div> <div>^"ANANAN"</div> </div>	
Add 7	Sort 7	Add 8	Sort 8	
<div> <div>^"BANANAN"</div> <div>^"BANANAN"</div> <div>^"BANANAN"</div> <div>^"BANANAN"</div> <div>^"BANANAN"</div> <div>^"BANANAN"</div> <div>^"BANANAN"</div> <div>^"BANANAN"</div> </div>	<div> <div>^"ANANANA"</div> <div>^"ANANANA"</div> <div>^"ANANANA"</div> <div>^"ANANANA"</div> <div>^"ANANANA"</div> <div>^"ANANANA"</div> <div>^"ANANANA"</div> <div>^"ANANANA"</div> </div>	<div> <div>^"BANANANA"</div> <div>^"BANANANA"</div> <div>^"BANANANA"</div> <div>^"BANANANA"</div> <div>^"BANANANA"</div> <div>^"BANANANA"</div> <div>^"BANANANA"</div> <div>^"BANANANA"</div> </div>	<div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> </div>	
Output				
<div> <div>^"BANANANA"</div> </div>				

Optimization

A number of optimizations can make these algorithms run more efficiently without changing the output. There is no need to represent the table in either the encoder or decoder. In the encoder, each row of the table can be represented by a single pointer into the strings, and the sort performed using the indices. Some care must be taken to ensure that the sort does not exhibit bad worst-case behavior: Standard library sort functions are unlikely to be appropriate. In the decoder, there is also no need to store the table, and in fact no sort is needed at all. In time proportional to the alphabet size and string length, the decoded string may be generated one character at a time from right to left. A "character" in the algorithm can be a byte, or a bit, or any other convenient size.

One may also make the observation that mathematically, the encoded string can be computed as a simple modification of the suffix array, and suffix arrays can be computed with linear time and memory. The BWT can be defined with regards to the suffix array SA of text T as (1-based indexing):

BWT[*i*] =

{

T
[
S
A
[
i
]
−
1
]
,

if
S
A
[
i
]
>
1

s
,

otherwise

{\displaystyle BWT[i]=\begin{cases}T[SA[i]-1], & {\text{if }}SA[i]>1\\s, & {\text{otherwise}}\end{cases}}

There is no need to have an actual 'EOF' character. Instead, a pointer can be used that remembers where in a string the 'EOF' would be if it existed. In this approach, the output of the BWT must include both the transformed string, and the final value of the pointer. That means the BWT does expand its input slightly. The inverse transform then shrinks it back down to the original size: it is given a string and a pointer, and returns just a string.

A complete description of the algorithms can be found in Burrows and Wheeler's paper, or in a number of online sources.

Bijective variant

When a bijective variant of the Burrows–Wheeler transform is performed on "BANANA", you get ANNBAA^ without the need for a special character for the end of the string. A special character forces one to increase character space by one, or to have a separate field with a numerical value for an offset. Either of these features makes data compression more difficult. When dealing with short files, the savings are great percentage-wise.

The bijective transform is done by sorting all rotations of the Lyndon words. In comparing two strings of unequal length, one can compare the infinite periodic repetitions of each of these in lexicographic order and take the last column of the base-rotated Lyndon word. For example, the text "BANANA" is transformed into "ANNBAA^" through these steps (the red ^ character indicates the EOF pointer) in the original string. The EOF character is unneeded in the bijective transform, so it is dropped during the transform and re-added to its proper place in the file.

The string is broken into Lyndon words so the words in the sequence are decreasing using the comparison method above. "BANANA" becomes (^) (B) (AN) (AN) (A), but Lyndon words are combined into (^) (B) (ANAN) (A).

Bijective transformation					
Input	All rotations	Sorted alphabetically by first letter	Last column of rotated Lyndon word	Output	
	<div> <div>^"BANANA"</div> <div>^"ANNBAAA"</div> <div>^"ANNBAAA"</div> <div>^"ANNBAAA"</div> <div>^"ANNBAAA"</div> <div>^"ANNBAAA"</div> <div>^"ANNBAAA"</div> <div>^"ANNBAAA"</div> </div>	<div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> </div>	<div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> </div>	<div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> <div>^"ANANANAN"</div> </div>	

Inverse bijective transform				
Input				
<div> <div>^"ANNBAA"</div> </div>				
Add 1	Sort 1	Add 2	Sort 2	
<div> <div>^"A"</div> <div>^"B"</div> <div>^"N"</div> <div>^"N"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> </div>	<div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> </div>	<div> <div>^"A"</div> <div>^"B"</div> <div>^"N"</div> <div>^"N"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> </div>	<div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> <div>^"A"</div> </div>	
Add 3	Sort 3	Add 4	Sort 4	
<div> <div>^"AN"</div> <div>^"AN"</div> <div>^"AN"</div> <div>^"AN"</div> <div>^"AN"</div> <div>^"AN"</div> <div>^"AN"</div> <div>^"AN"</div> </div>	<div> <div>^"ANA"</div> <div>^"ANA"</div> <div>^"ANA"</div> <div>^"ANA"</div> <div>^"ANA"</div> <div>^"ANA"</div> <div>^"ANA"</div> <div>^"ANA"</div> </div>	<div> <div>^"ANAN"</div> <div>^"ANAN"</div> <div>^"ANAN"</div> <div>^"ANAN"</div> <div>^"ANAN"</div> <div>^"ANAN"</div> <div>^"ANAN"</div> <div>^"ANAN"</div> </div>	<div> <div>^"ANANAN"</div> <div>^"ANANAN"</div> <div>^"ANANAN"</div> <div>^"ANANAN"</div> <div>^"ANANAN"</div> <div>^"ANANAN"</div> <div>^"ANANAN"</div> <div>^"ANANAN"</div> </div>	
Output				
<div> <div>^"BANANANA"</div> </div>				

The above may be viewed as four cycles:

^ = (^)(^)... = ^^^^^...

B = (B)(B)... = BBBB...

ANAN = (AN)(AN)... = ANANANAN...

A = (A)(A)... = AAAAA...

or 5 cycles WHERE ANAN broken into 2:

AN = (AN) (AN) ... = ANANANAN

AN = (AN) (AN) ... = ANANANAN

If a cycle is N character it will be repeated N times:

(^)

(B)

(ANAN)

(A)

or

(^)

(B)

(AN)

(AN)

(A)

to get the

^BANANANA

Since any rotation of the input string will lead to the same transformed string, the BWT cannot be inverted without adding an EOF marker to the input or, augmenting the output with information such as an index, making it possible to identify the input string from all its rotations.

There is a bijective version of the transform, by which the transformed string uniquely identifies the original. In this version, every string has a unique inverse of the same length.^{[4][5]}

The fastest versions are linear in time and space.

The bijective transform is computed by factoring the input into a non-increasing sequence of Lyndon words; such a factorization exists in the Chen–Fox–Lyndon theorem,^[6] and may be found in linear time.^[7] The algorithm sorts the rotations of all the words; as in the Burrows–Wheeler transform, this produces a sorted sequence of *n* strings. The transformed string is then obtained by picking the final character of each string in this sorted list.

For example, applying the bijective transform gives:

Input	SIX.MIXED.PIXIES.SIFT.SIXTY.PIXIE.DUST.BOXES
Lyndon words	SIX.MIXED.PIXIES.SIFT.SIXTY.PIXIE.DUST.BOXES
Output	STEYDST.E.IXIXIXXSMPPFS.B..EE..SUSFXDIOIIIT

The bijective transform includes eight runs of identical characters. These runs are, in order: XX, II, XX, PP, . . . , EE, . . . , and IIII.

In total, 18 characters are used in these runs.

Dynamic Burrows–Wheeler transform

When a text is edited, its Burrows–Wheeler transform changes. Salson *et al.*^[8] propose an algorithm that deduces the Burrows–Wheeler transform of an edited text from that of the original text, doing a limited number of local re-orderings in the original Burrows–Wheeler transform, which can be faster than constructing the Burrows–Wheeler transform of the edited text directly.

Sample implementation

This Python implementation sacrifices speed for simplicity: the program is short, but takes more than the linear time that would be desired in a practical implementation.

Using the STX/ETX control codes to mark the start and end of the text, and using s[i:] + s[:i] to construct the ith rotation of s, the forward transform takes the last character of each of the sorted rows:

```

def bwt(s):
    """Apply Burrows-Wheeler transform to input string.***

```