

Burrows–Wheeler transform

The **Burrows–Wheeler transform** (**BWT**, also called **block-sorting compression**), is an algorithm used in data compression techniques such as bzip2. It was invented by Michael Burrows and David Wheeler in 1994 while working at DEC Systems Research Center in Palo Alto, California.^[1] It is based on a previously unpublished transformation discovered by Wheeler in 1983.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, none of its characters change value. 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. 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.

For example:

Input	SIX.MIXED.PIXIES.SIFT.SIXTY.PIXIE.DUST.BOXES
Output	TEXYDST.E.IXIXIXSSMPPS.B..E.S.EUSFXDIIIOIIIT

The output is easier to compress because it has many repeated characters. In fact, in the transformed string, there are a total of six runs of identical characters: XX, SS, PP, . . , II, and III, which together make 13 out of the 44 characters in it.

Example

The transform is done by sorting all rotations of the text in lexicographic order, then taking the last column. 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 List of Rotations	Output Last Column
^BANANA@	^BANANA@	ANANA@ ^B	BNN^AA@A
	@^BANANA	ANA@ ^BAN	
	A@^BANAN	A@ ^BANAN	
	NA@^BANA	BANANA@ ^	
	ANA@^BAN	NANA@ ^BA	
	NANA@^BA	NA@ ^BANA	
	ANANA@^B	^BANANA@	
	BANANA@^	@ ^BANANA	

The following pseudocode gives a simple but 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, and is ignored during sorting.

```
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
    insert s as a column of table before first column of the table // first insert creates first column
    sort rows of the table alphabetically
return (row that ends with the 'EOF' character)
```

To understand why this creates more-easily-compressible data, let's consider transforming a long English text frequently containing the word "the". Sorting the rotations of this text will often group rotations starting with "he " 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 "Brahe ") 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 to get the first column. Then, the first and last columns 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			
BNN^AA@A			
Add 1	Sort 1	Add 2	Sort 2
B	A	BA	AN
N	A	NA	AN
N	A	NA	A@
^	B	^B	BA
A	N	AN	NA
A	N	AN	NA
@	^	@^	^B
A	@	A@	@^
Add 3	Sort 3	Add 4	Sort 4
BAN	ANA	BANA	ANAN
NAN	ANA	NANA	ANA@
NA@	A@^	NA@^	A@^B
^BA	BAN	^BAN	BANA
ANA	NAN	ANAN	NANA
ANA	NA@	ANA@	NA@^
@^B	^BA	@^BA	^BAN
A@^	@^B	A@^B	@^BA
Add 5	Sort 5	Add 6	Sort 6

BANAN	ANANA	BANANA	ANANA
NANA	ANA	NANA	ANA
NA	A	NA	A
^BANA	BANAN	^BANAN	BANANA
ANANA	NANA	ANANA	NANA
ANA	NA	ANA	NA
^BAN	^BANA	^BANANA	^BANAN
A	^BAN	A	^BAN
Add 7	Sort 7	Add 8	Sort 8
BANANA	ANANA	BANANA	ANANA
NANA	ANA	NANA	ANA
NA	A	NA	A
^BANANA	BANANA	^BANANA	BANANA
ANANA	NANA	ANANA	NANA
ANA	NA	ANA	NA
^BANAN	^BANANA	^BANANA	^BANANA
A	^BANAN	A	^BANAN
Output			
^BANANA			

Optimization

A number of optimizations can make these algorithms run more efficiently without changing the output. In BWT, 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.

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.

Bijjective variant

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, that makes it possible to identify the input string from the class of all of 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.^[2]

The bijective transform is computed by first factoring the input into a non-increasing sequence of Lyndon words; such a factorization exists by the Chen–Fox–Lyndon theorem, and can be found in linear time.^[3] Then, the algorithm sorts together all the rotations of all of these words; as in the usual Burrows-Wheeler transform, this produces a sorted sequence of n strings. The transformed string is then obtained by picking the final character of each of these strings 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.IXXIIXXSMPPXS.B..EE..SUSFXDIOIIIIIT

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 take part in these runs.

Dynamic Burrows–Wheeler transform

Instead of reconstructing the Burrows–Wheeler transform of an edited text, Salson *et al.*^[4] propose an algorithm that deduces the new Burrows–Wheeler transform from the original one, doing a limited number of local reorderings in the original Burrows–Wheeler transform.

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 null character as the end of file marker, and using `s[i:] + s[:i]` to construct the *i*th 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."""
    assert "\0" not in s, "Input string cannot contain null character"
    ('\0')
    s += "\0" # Add end of file marker
    table = sorted(s[i:] + s[:i] for i in range(len(s))) # Table of
rotations of string
    last_column = [row[-1:] for row in table] # Last characters of
each row
    return "".join(last_column) # Convert list of characters into
string
```

The inverse transform repeatedly inserts `r` as the left column of the table and sorts the table. After the whole table is built, it returns the row that ends with null, minus the null.

```
def ibwt(r):
    """Apply inverse Burrows–Wheeler transform."""
    table = [""] * len(r) # Make empty table
    for i in range(len(r)):
        table = sorted(r[i] + table[i] for i in range(len(r))) # Add a
column of r
    s = [row for row in table if row.endswith("\0")][0] # Find the
correct row (ending in "\0")
    return s.rstrip("\0") # Get rid of trailing null character
```

BWT in bioinformatics

The advent of high-throughput sequencing (HTS) techniques at the end of the 2000 decade has led to another application of the Burrows–Wheeler transformation. In HTS, DNA is fragmented into small pieces, of which the first few bases are sequenced, yielding several millions of "reads", each 30 to 500 base pairs ("DNA characters") long. In many experiments, e.g., in ChIP-Seq, the task is now to align these reads to a reference genome, i.e., to the known, nearly complete sequence of the organism in question (which may be up to several billion base pairs long). A number of alignment programs, specialized for this task, were published, which initially relied on hashing (e.g., Eland, SOAP^[5],^[6] or Maq^[7]). In an effort to reduce the memory requirement for sequence alignment, several alignment programs were developed (Bowtie,^[8] BWA,^[9] and SOAP2^[10]) which use the Burrows–Wheeler transform.

References

- [1] Burrows M and Wheeler D (1994), *A block sorting lossless data compression algorithm* (<http://www.hpl.hp.com/techreports/Compaq-DEC/SRC-RR-124.html>), Technical Report 124, Digital Equipment Corporation,
- [2] Gil, J.; Scott, D. A. (2009), *A bijective string sorting transform* (<http://bijective.dogma.net/00yyy.pdf>), . Kufleitner, Manfred (2009), "On bijective variants of the Burrows–Wheeler transform" (<http://www.stringology.org/event/2009/p07.html>), in Holub, Jan; Žďárek, Jan, *Prague Stringology Conference*, pp. 65–69, arXiv:0908.0239, .
- [3] Duval, Jean-Pierre (1983), "Factorizing words over an ordered alphabet", *Journal of Algorithms* **4** (4): 363–381, doi:10.1016/0196-6774(83)90017-2.
- [4] Salson M, Lecroq T, Léonard M and Mouchard L (2009). "A Four-Stage Algorithm for Updating a Burrows–Wheeler Transform". *Theoretical Computer Science* **410** (43): 4350. doi:10.1016/j.tcs.2009.07.016.
- [5] <http://soap.genomics.org.cn>
- [6] Li R, *et al.* (2008). "SOAP: short oligonucleotide alignment program". *Bioinformatics* **24** (5): 713–714. doi:10.1093/bioinformatics/btn025. PMID 18227114.
- [7] Li H, Ruan J, Durbin R (2008-08-19). "Mapping short DNA sequencing reads and calling variants using mapping quality scores". *Genome Research* **18** (11): 1851–1858. doi:10.1101/gr.078212.108. PMC 2577856. PMID 18714091.
- [8] Langmead B, Trapnell C, Pop M, Salzberg SL (2009). "Ultrafast and memory-efficient alignment of short DNA sequences to the human genome". *Genome Biology* **10** (3): R25. doi:10.1186/gb-2009-10-3-r25. PMC 2690996. PMID 19261174.
- [9] Li H, Durbin R (2009). "Fast and accurate short read alignment with Burrows–Wheeler Transform". *Bioinformatics* **25** (14): 1754–1760. doi:10.1093/bioinformatics/btp324. PMC 2705234. PMID 19451168.
- [10] Li R, *et al.* (2009). "SOAP2: an improved ultrafast tool for short read alignment". *Bioinformatics* **25** (15): 1966–1967. doi:10.1093/bioinformatics/btp336. PMID 19497933.

External links

- Compression comparison of BWT based file compressors (<http://compressionratings.com/bwt.html>)
- Article by Mark Nelson on the BWT (<http://marknelson.us/1996/09/01/bwt/>)
- A Bijective String-Sorting Transform, by Gil and Scott (<http://bijective.dogma.net/00yyy.pdf>)
- On Bijective Variants of the Burrows–Wheeler Transform, by Kufleitner (<http://arxiv.org/abs/0908.0239>)
- Blog post (<http://google-opensource.blogspot.com/2008/06/debuting-dcs-bwt-experimental-burrows.html>) and project page (<http://code.google.com/p/dcs-bwt-compressor/>) for an open-source compression program and library based on the Burrows–Wheeler algorithm

Article Sources and Contributors

Burrows–Wheeler transform *Source:* <http://en.wikipedia.org/w/index.php?oldid=457617624> *Contributors:* 130.94.122.xxx, 213.253.39.xxx, Alexr, Ambulnick, Beland, Bloodhold, Brighterorange, Brion VIBBER, Burton Radons, Cbogart2, Charles Matthews, Connelly, Cyhawk, Cyp, Damian Yerrick, DataWraith, David Eppstein, Dcoetzee, Dicklyon, DmitriyV, Doradus, Drachefly, Drnathanfurious, Ed g2s, Faisal.akeel, Felix Wiemann, Fredrik, Garandel, Giftlite, GregorB, Gstein, Henning Makholm, Inkling, Intgr, JakeVortex, Jaredwf, Jerry, John K, Kku, LC, Malbrain, Mark.t.nelson, MarkHudson, Mc6809e, Michael Hardy, Mormegil, Nikai, Ocolon, Okted, Oli Filth, OverlordQ, PMLawrence, PierreAbbat, Piet Delpont, Pne, Populus, Pt, Quuxplusone, R. S. Shaw, Reedy, Requestion, Robackja, RolandIllig, Rparle, Rursus, Samir000, Saxbryn, Sligocki, Speight, Spooky, Stw, Taak, Taw, TeeEmCee, Thesuperslacker, Thorwald, Timwi, Torsten Will, Tribaal, Wfaulk, WhiteDragon, WikiReviewer.de, Wikibofh, Xchmelmilos, Xpicto, Yanghoch, ZeroOne, 90 anonymous edits

License

Creative Commons Attribution-Share Alike 3.0 Unported
[//creativecommons.org/licenses/by-sa/3.0/](http://creativecommons.org/licenses/by-sa/3.0/)