# Burrows-Wheeler Transform, part 1

Ben Langmead

JOHNS HOPKINS
WHITING SCHOOL
*of* ENGINEERING
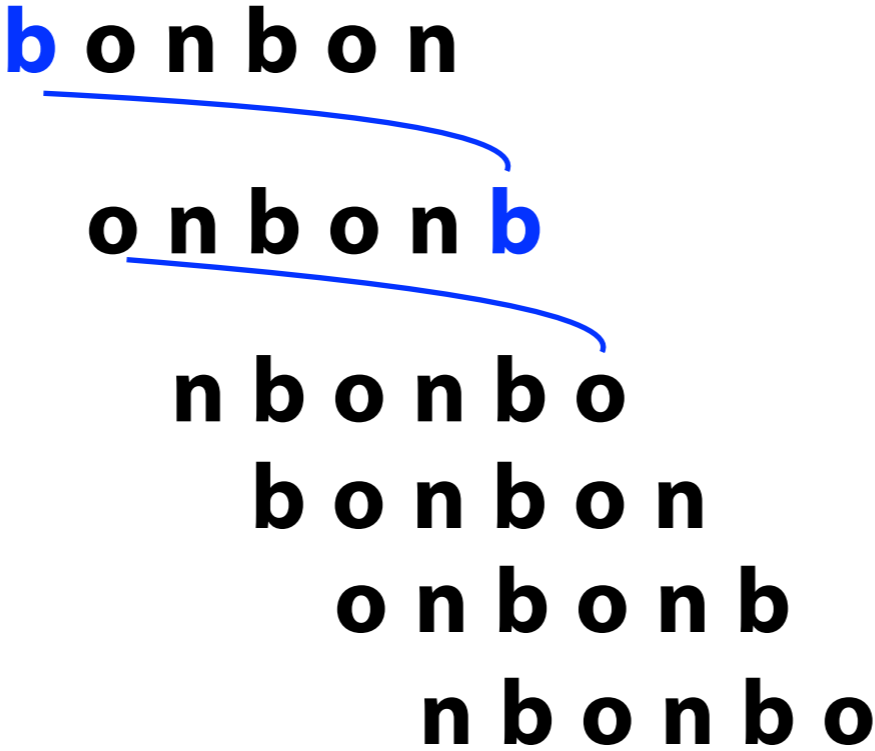
## Department of Computer Science

# Burrows-Wheeler Transform

Rotations of a string:

**b** o n b o n

o n b o n **b**

**n b o n b o**
**b o n b o n**
**o n b o n b**
**n b o n b o**

(after this they
repeat)

bonbon

nbonbo

onbonb

bonbon

nbonbo

onbonb

# Burrows-Wheeler Transform

We know dictionary order:
`as < ash` and `flower < flowers`

This is a *convention* for cases where no character comparison "breaks the tie," i.e. where one string is a prefix of the other

We could have said `ash < as` and
`flowers < flower;` still a total order

# Burrows-Wheeler Transform

We invent a new symbol **$** ("terminator"), defined to be alphabetically less than all others:

bonbon**$**

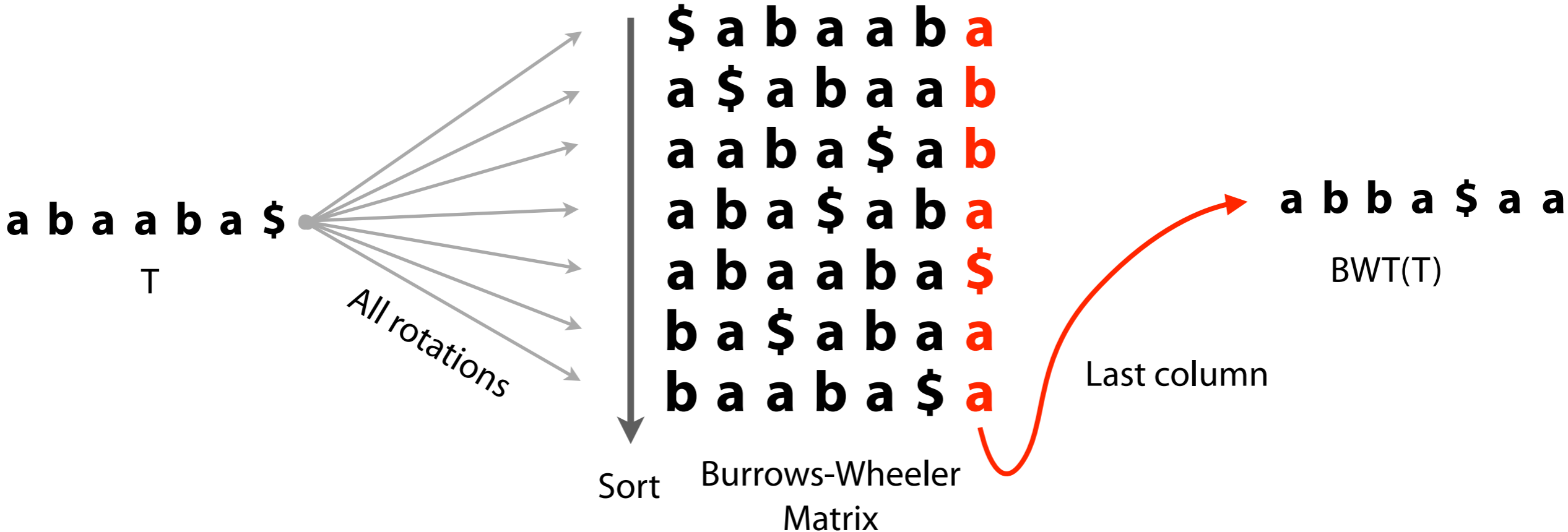Enforces dictionary order and:

No suffix is a prefix of another suffix:

bonbon**$**
onbon**$**
nbon**$**
bon**$**
on**$**
n**$**
**$**

And no two rotations are the same:

bonbon**$**
**$**bonbon
n**$**bonbo
on**$**bonb
bon**$**bon
nbon**$**bo
onbon**$**b

# Burrows-Wheeler Transform



a b a a b a $

T

All rotations

Sort

$ a b a a b **a**
a $ a b a a **b**
a a b a $ a **b**
a b a $ a b **a**
a b a a b a **$**
b a $ a b a **a**
b a a b a $ **a**

Burrows-Wheeler
Matrix

Last column

a b b a $ a a

BWT(T)

Burrows M, Wheeler DJ: A block sorting lossless data compression algorithm.
*Digital Equipment Corporation, Palo Alto, CA* 1994, Technical Report 124; 1994

# Burrows-Wheeler Transform

**a b a a b a $**

T

| $ a b a a b **a** | | **a** $ a b a a b |
|---|---|---|
| a $ a b a a **b** | | **b** a $ a b a a |
| a a b a $ a **b** | | **b** a a b a $ a |
| a b a $ a b **a** | → | **a** a b a $ a b |
| a b a a b a **$** | | **$** a b a a b a |
| b a $ a b a **a** | | **a** b a $ a b a |
| b a a b a $ **a** | | **a** b a a b a $ |

BWT    Right contexts

BWT(T) orders T's characters according to alphabetical order of right contexts in T

# Burrows-Wheeler Transform

Ordered by *right-context*

Colors show what parts of
matrix are shown on right

| final char (L) | sorted rotations |
|---|---|
| a | n to decompress.   It achieves compression |
| o | n to perform only comparisons to a depth |
| o | n transformation}   This section describes |
| o | n transformation}   We use the example and |
| o | n treats the right-hand side as the most |
| a | n tree for each 16 kbyte input block, enc |
| a | n tree in the output stream, then encodes |
| i | n turn, set $L[i]$ to be the |
| i | n turn, set $R[i]$ to the |
| o | n unusual data. Like the algorithm of Man |
| a | n use a single set of probabilities table |
| e | n using the positions of the suffixes in |
| i | n value at a given point in the vector $R |
| e | n we present modifications that improve t |
| e | n when the block size is quite large.  Ho |
| i | n which codes that have not been seen in |
| i | n with $ch$ appear in the {\em same order |
| i | n with $ch$.                    In our exam |
| o | n with Huffman or arithmetic coding.  Bri |
| o | n with figures given by Bell~\cite{bell}. |

Figure 1: Example of sorted rotations. Twenty consecutive rotations from the sorted list of rotations of a version of this paper are shown, together with the final character of each rotation.

Figure: Burrows M, Wheeler DJ: A block sorting lossless data compression algorithm.
*Digital Equipment Corporation, Palo Alto, CA* 1994, Technical Report 124; 1994

# Burrows-Wheeler Transform

Consider building a
high-order compressor
for `mississippi`

```
$ m i s s i s s i p p i
i $ m i s s i s s i p p
i p p i $ m i s s i s s
i s s i p p i $ m i s s
i s s i s s i p p i $ m
m i s s i s s i p p i $
p i $ m i s s i s s i p
p p i $ m i s s i s s i
s i p p i $ m i s s i s
s i s s i p p i $ m i s
s s i p p i $ m i s s i
s s i s s i p p i $ m i
```

Say we're using **right**
context, $k = 1$

$$H_1(T) = (1/12)\, H_0(\texttt{i}) + (4/12)\, H_0(\texttt{pssm}) +$$
$$= (1/12)\, H_0(\texttt{\$}) + (2/12)\, H_0(\texttt{pi}) + (4/12)\, H_0(\texttt{ssii})$$

# Burrows-Wheeler Transform

| | |
|---|---|
| $\$$mississippi | $H_0$ |
| i$\$$mississipp | |
| ippi$\$$mississ | |
| issippi$\$$miss | $H_0$ |
| ississippi$\$$m | |
| mississippi$\$$ | $H_0$ |
| pi$\$$mississip | |
| ppi$\$$mississi | $H_0$ |
| sippi$\$$missis | |
| sissippi$\$$mis | |
| ssippi$\$$missi | $H_0$ |
| ssissippi$\$$mi | |

Overall: $H_1$

We obtain a $H_k$ compressor by using a $H_0$ compressor in each length-$k$-context chunk. $k$ can vary.

# Burrows-Wheeler Transform

| | |
|---|---|
| $\$$ m i s s i s s i p p i | $H_0$ |
| i $\$$ m i s s i s s i p p | $H_0$ |
| i p p i $\$$ m i s s i s s | $H_0$ |
| i s s i p p i $\$$ m i s s | |
| i s s i s s i p p i $\$$ m | $H_0$ |
| m i s s i s s i p p i $\$$ | $H_0$ |
| p i $\$$ m i s s i s s i p | $H_0$ |
| p p i $\$$ m i s s i s s i | $H_0$ |
| s i p p i $\$$ m i s s i s | |
| s i s s i p p i $\$$ m i s | $H_0$ |
| s s i p p i $\$$ m i s s i | |
| s s i s s i p p i $\$$ m i | $H_0$ |

Overall: $H_2$

We obtain a $H_k$ compressor by using a $H_0$ compressor in each length-$k$-context chunk. $k$ can vary.

# Burrows-Wheeler Transform

Compression strategy:

   **(a)** Find BWT(T)

   **(b)** Partition BWT by $k$-context

   **(c)** Apply $H_0$ encoding in each

Space required:

  $H_0$ codebook for
each partition

bzip2 uses this method

Glossing over
details here

Decompression strategy:

   **(a)** Decode $H_0$ in each partition

   **(b)** Concatenate partitions

   **(c)** Find T by reversing BWT(T)

TODO

BWT is a "compression booster"