# Reversing the BWT

Ben Langmead

JOHNS HOPKINS
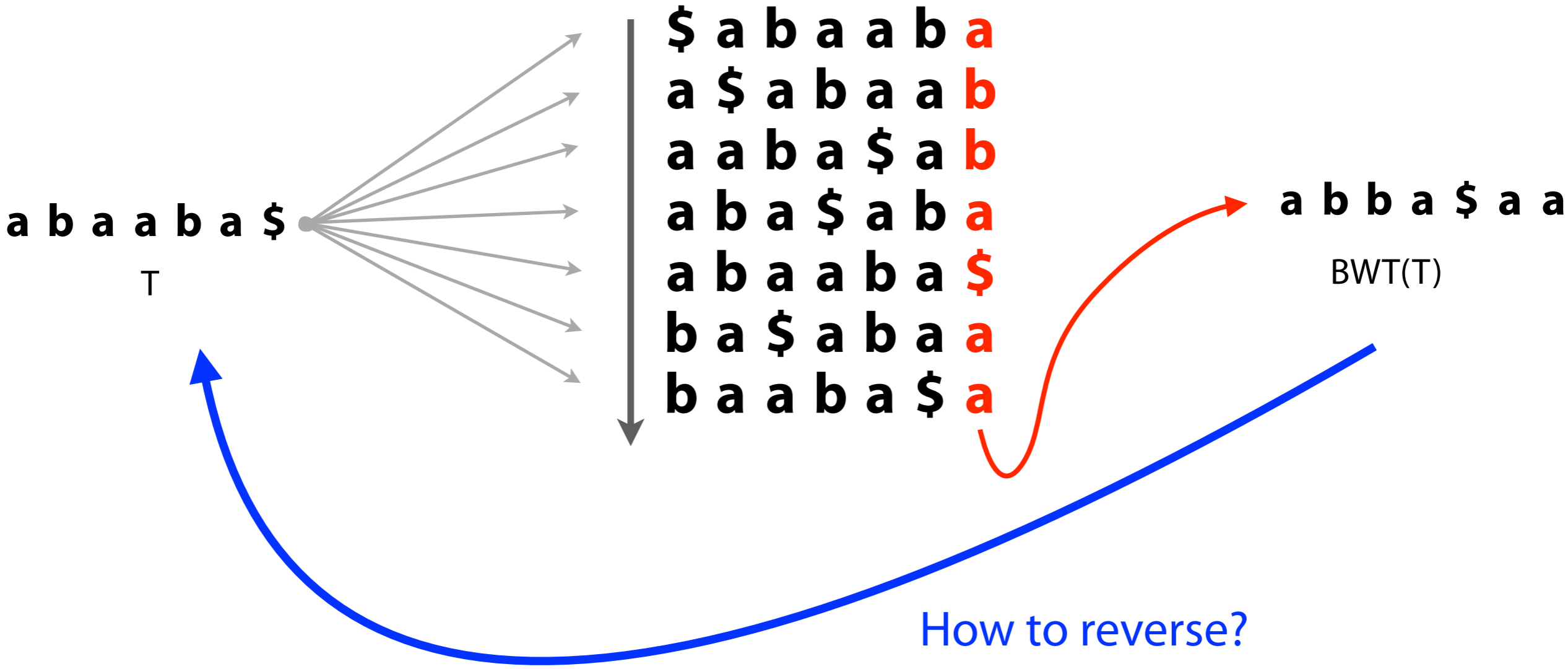WHITING SCHOOL
*of* ENGINEERING

## Department of Computer Science

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

a b b a $ a a
BWT(T)

How to reverse?

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

Give each character in *T* a ***rank:***

**a   b   a   a   b   a   $**

# Burrows-Wheeler Transform

Give each character in *T* a ***rank:***

$$a_0 \ b_0 \ a_1 \ a_2 \ b_1 \ a_3 \ \$$$

Ranks aren't explicitly stored; we use them to distinguish occurrences

Let's re-write the BWM with ranks...

# Burrows-Wheeler Transform

BWM with ranks:

|  | $F$ |  |  |  |  |  | $L$ |
|---|---|---|---|---|---|---|---|
| $\$$ | $a_0$ | $b_0$ | $a_1$ | $a_2$ | $b_1$ | $a_3$ | |
| $a_3$ | $\$$ | $a_0$ | $b_0$ | $a_1$ | $a_2$ | $b_1$ | |
| $a_1$ | $a_2$ | $b_1$ | $a_3$ | $\$$ | $a_0$ | $b_0$ | |
| $a_2$ | $b_1$ | $a_3$ | $\$$ | $a_0$ | $b_0$ | $a_1$ | |
| $a_0$ | $b_0$ | $a_1$ | $a_2$ | $b_1$ | $a_3$ | $\$$ | |
| $b_1$ | $a_3$ | $\$$ | $a_0$ | $b_0$ | $a_1$ | $a_2$ | |
| $b_0$ | $a_1$ | $a_2$ | $b_1$ | $a_3$ | $\$$ | $a_0$ | |

Look at first and last columns, called $F$ and $L$

And look at just the **a**s

**a**s occur in the same order in $F$ and $L$.  As we look down columns, in both cases we see:  **a**$_3$, **a**$_1$, **a**$_2$, **a**$_0$

# Burrows-Wheeler Transform

BWM with ranks:

|  |  |  | $F$ |  |  |  |  |  |  |  | $L$ |
|---|---|---|---|---|---|---|---|---|---|---|---|

$$\$ \quad a_0 \quad b_0 \quad a_1 \quad a_2 \quad b_1 \quad a_3$$

$$a_3 \quad \$ \quad a_0 \quad b_0 \quad a_1 \quad a_2 \quad \mathbf{b_1}$$

$$a_1 \quad a_2 \quad b_1 \quad a_3 \quad \$ \quad a_0 \quad \mathbf{b_0}$$

$$a_2 \quad b_1 \quad a_3 \quad \$ \quad a_0 \quad b_0 \quad a_1$$

$$a_0 \quad b_0 \quad a_1 \quad a_2 \quad b_1 \quad a_3 \quad \$$$

$$\mathbf{b_1} \quad a_3 \quad \$ \quad a_0 \quad b_0 \quad a_1 \quad a_2$$

$$\mathbf{b_0} \quad a_1 \quad a_2 \quad b_1 \quad a_3 \quad \$ \quad a_0$$

Same with **b**s:   $\mathbf{b_1}$, $\mathbf{b_0}$

# Burrows-Wheeler Transform

BWM with ranks:

| $F$ | | | | | | $L$ |
|---|---|---|---|---|---|---|
| $ | $a_0$ | $b_0$ | $a_1$ | $a_2$ | $b_1$ | $a_3$ |
| $a_3$ | $ | $a_0$ | $b_0$ | $a_1$ | $a_2$ | $b_1$ |
| $a_1$ | $a_2$ | $b_1$ | $a_0$ | $ | $a_0$ | $b_0$ |
| $a_2$ | $b_1$ | $a_3$ | $ | $a_0$ | $b_0$ | $a_1$ |
| $a_0$ | $b_0$ | $a_1$ | $a_2$ | $b_1$ | $a_3$ | $ |
| $b_1$ | $a_3$ | $ | $a_0$ | $b_0$ | $a_1$ | $a_2$ |
| $b_0$ | $a_1$ | $a_2$ | $b_1$ | $a_3$ | $ | $a_0$ |

LF Mapping: The $i^{th}$ occurrence of a character $c \in \Sigma$ in $L$ and the $i^{th}$ occurrence of $c$ in $F$ correspond to the *same* occurrence in $T$ (i.e. have ***same rank***)

# Burrows-Wheeler Transform

Why does the LF Mapping hold?

Why are these **a**s in this order relative to each other?

```
$  a  b  a  a  b  a₃
a₃ $  a  b  a  a  b₁
a₁ a  b  a  $  a  b₀
a₂ b  a  $  a  b  a₁
a₀ b  a  a  b  a  $
b₁ a  $  a  b  a  a₂
b₀ a  a  b  a  $  a₀
```

They're sorted by right-context

Why are these **a**s in this order relative to each other?

```
$  a  b  a  a  b  a₃
a₃ $  a  b  a  a  b₁
a₁ a  b  a  $  a  b₀
a₂ b  a  $  a  b  a₁
a₀ b  a  a  b  a  $
b₁ a  $  a  b  a  a₂
b₀ a  a  b  a  $  a₀
```

They're sorted by right-context

Occurrences of $c$ in $F$ are sorted by right-context; same for $L$

# Burrows-Wheeler Transform

Reverse BWT(T) starting at right end of *T*, moving left

**Start** in first row. *F* must have **$**.

*L* contains character prior: **a$_3$**

**Jump** to row *beginning* with **a$_3$**.

*L* contains character just prior: **b$_1$**.

Repeat for **b$_1$**, get **a$_2$**

Repeat for **a$_2$**, get **a$_1$**

Repeat for **a$_1$**, get **b$_0$**

Repeat for **b$_0$**, get **a$_0$**

Repeat for **a$_0$**, get **$** (done)

*T:* **a$_0$ b$_0$ a$_1$ a$_2$ b$_1$ a$_3$ $**

*F*          *L*

**$**  ⟶  **a$_3$**

**a$_3$**  ⟶  **b$_1$**

**a$_1$**  ⟶  **b$_0$**

**a$_2$**  ⟶  **a$_1$**

**a$_0$**  ⟶  **$**

**b$_1$**  ⟶  **a$_2$**

**b$_0$**  ⟶  **a$_0$**

We visited (backwards) T's chars:

**a$_0$ b$_0$ a$_1$ a$_2$ b$_1$ a$_3$ $**