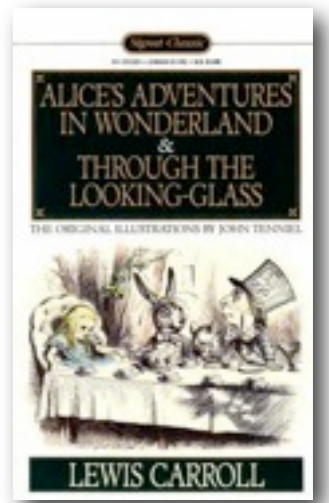# Burrows-Wheeler Transform

## CMSC 423

# Genome of the Cow

a sequence of 2.86 billion letters

enough letters to fill a million pages of a typical book.

```
TATGGAGCCAGGTGCCTGGGGCAACAAGACTGTGGTCACTGAATTCATCCTTCTTGGTCTAACAGAGAACATAG
AACTGCAATCCATCCTTTTTGCCATCTTCCTCTTTGCCTATGTGATCACAGTCGGGGGCAACTTGAGTATCCTG
GCCGCCATCTTTGTGGAGCCCAAACTCCACACCCCCATGTACTACTTCCTGGGGAACCTTTCTCTGCTGGACAT
TGGGTGCATCACTGTCACCATTCCTCCCATGCTGGCCTGTCTCCTGACCCACCAATGCCGGGTTCCCTATGCAG
CCTGCATCTCACAGCTCTTCTTTTTCCACCTCCTGGCTGGAGTGGACTGTCACCTCCTGACAGCCATGGCCTAC
GACCGCTACCTGGCCATTTGCCAGCCCCTCACCTATAGCATCCGCATGAGCCGTGACGTCCAGGGAGCCCTGGT
GGCCGTCTGCTGCTCCATCTCCTTCATCAATGCTCTGACCCACACAGTGGCTGTGTCTGTGCTGGACTTCTGCG
GCCCTAACGTGGTCAACCACTTCTACTGTGACCTCCCGCCCCTTTTCCAGCTCTCCTGCTCCAGCATCCACCTC
AACGGGCAGCTACTTTTCGTGGGGGCCACCTTCATGGGGGTGGTCCCCATGGTCTTCATCTCGGTATCCTATGC
CCACGTGGCAGCCGCAGTCCTGCGGATCCGCTCGGCAGAGGGCAGGAAGAAAGCCTTCTCCACGTGTGGCTCCC
ACCTCACCGTGGTCTGCATCTTTTATGGAACCGGCTTCTTCAGCTACATGCGCCTGGGCTCCGTGTCCGCCTCA
GACAAGGACAAGGGCATTGGCATCCTCAACACTGTCATCAGCCCCATGCTGAACCCACTCATCTACAGCCTCCG
GAACCCTGATGTGCAGGGCGCCCTGAAGAGGTTGCTGACAGGGAAGCGGCCCCCGGAGTG  . . .
```

# We can only read ~ 1000 characters at a time from a random place:

good-natured, she thought: still

          when it saw Alice. It looked

    ought to be treated

          good-natured, she thought, still

      Cat only

     a greet many

            It looked good-

The Cat only grinned when it saw Alice.

   be treated with respect.

       still it had very long claws

    claws and a great many teeth, so she

         so she felt that it ought

The Cat only grinned when it saw Alice.
    Cat only          when it saw Alice. It looked
                                    It looked good-

good-natured, she thought: still
good-natured, she thought, still
                        still it had very long claws

claws and a great many teeth, so she
        a greet many          so she felt that it ought

ought to be treated
        be treated with respect.



It's a jigsaw
puzzle ...

...except with 35
million pieces

# Motivation - Short Read Mapping

A Cow Genome

Sequencing technologies produce millions of "reads" = a random, short substring of the genome

If we already know the genome of one cow, we can get reads from a 2nd cow and map them onto the known cow genome.

Need to do millions of string searches in a long string.

# Burrows-Wheeler Transform

Text transform that is useful for compression & search.

banana

| banana$ | | $banana |
|---------|------|---------|
| anana$b | | a$banan |
| nana$ba | sort → | ana$ban |
| ana$ban | | anana$b |
| na$bana | | banana$ |
| a$banan | | nana$ba |
| $banana | | na$bana |

BWT(banana) = annb$aa

Tends to put runs of the same character together.

Makes compression work well.

"bzip" is based on this.

# Another Example

appellee$

appellee$
ppellee$a
pellee$ap
ellee$app
llee$appe    sort →
lee$appel
ee$appell
e$appelle
$appellee

$appellee
appellee$
e$appelle
ee$appell
ellee$app
lee$appel
llee$appe
pellee$ap
ppellee$a

BWT(appellee$) =
e$elplepa

Doesn't always improve
the compressibility...

# Recovering the string

BWT

sort → first column
BWT

→ first 2 columns

↗ first 3 columns

| | | |
|---|---|---|
| $appellee | e | |
| appellee$ | $ | |
| e$appelle | e | |
| ee$appell | l | |
| ellee$app | p | |
| lee$appel | l | |
| llee$appe | e | |
| pellee$ap | p | |
| ppellee$a | a | |

e  $
$  a
e  e
l  e
p  e
l  l
e  l
p  p
a  p

— sort these 2 columns →

$a
ap
e$
ee
el
le
ll
pe
pp

— prepend BWT column →

e  $a
$  ap
e  e$
l  ee
p  el
l  le
e  ll
p  pe
a  pp

— Sort these 3 columns →

$ap
app
e$a
ee$
ell
lee
lle
pel
ppe

# Inverse BWT

```
def inverseBWT(s):
    B = [s_1,s_2,s_3,...,s_n]
    for i = 1..n:
        sort B
        prepend s_i to B[i]
    return row of B that ends with $
```

# Another BWT Example

dogwood$
ogwood$d
gwood$do
wood$dog          sort          $dogwood
ood$dogw        $\longrightarrow$        d$dogwoo
od$dogwo                        dogwood$
d$dogwoo                        gwood$do          last column
$dogwood                        od$dogwo        $\longrightarrow$
                                ogwood$d          BWT(dogwood$) =
                                ood$dogw                do$oodwg
                                wood$dog

# Another BWT Example

**Prepend** | **Sort**
--- | ---
d $ | $d
o d | d$
$ d | do
o g | gw
o o | od
d o | og
w o | oo
g w | wo

**Prepend** | **Sort**
--- | ---
d $d | $do
o d$ | d$d
$ do | dog
o gw | gwo
o od | od$
d og | ogw
w oo | ood
g wo | woo

**Prepend** | **Sort**
--- | ---
d$do | $dog
od$d | d$do
$dog | dogw
ogwo | gwoo
ood$ | od$d
dogw | ogwo
wood | ood$
gwoo | wood

**Prepend** | **Sort**
--- | ---
d $dog | $dogw
o d$do | d$dog
$ dogw | dogwo
o gwoo | gwood
o od$d | od$do
d ogwo | ogwoo
w ood$ | ood$d
g wood | wood$

**Prepend** | **Sort**
--- | ---
d $dogw | $dogwo
o d$dog | d$dogw
$ dogwo | dogwoo
o gwood | gwood$
o od$do | od$dog
d ogwoo | ogwood
w ood$d | ood$do
g wood$ | wood$d

**Prepend** | **Sort**
--- | ---
d $dogwo | $dogwoo
o d$dogw | d$dogwo
$ dogwoo | dogwood
o gwood$ | gwood$d
o od$dog | od$dogw
d ogwood | ogwood$
w ood$do | ood$dog
g wood$d | wood$do

**Prepend** | **Sort**
--- | ---
d $dogwoo | $dogwood
o d$dogwo | d$dogwoo
$ dogwood | dogwood$
o gwood$d | gwood$do
o od$dogw | od$dogwo
d ogwood$ | ogwood$d
w ood$dog | ood$dogw
g wood$do | wood$dog

# Searching with BWT: LF Mapping

LF Mapping

BWT(unabashable)

| | $ | a | b | e | h | l | n | s | u |
|---|---|---|---|---|---|---|---|---|---|
| $unabashable | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| abashable$un | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| able$unabash | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| ashable$unab | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| bashable$una | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| ble$unabasha | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| e$unabashabl | 0 | 2 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| hable$unabas | 0 | 2 | 1 | 1 | 1 | 1 | 1 | 0 | 0 |
| le$unabashab | 0 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 0 |
| nabashable$u | 0 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 0 |
| shable$unaba | 0 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 1 |
| unabashable$ | 0 | 3 | 2 | 1 | 1 | 1 | 1 | 1 | 1 |
| | 1 | 3 | 2 | 1 | 1 | 1 | 1 | 1 | 1 |

# of times letter appears before this position in the last column.

**LF Property:** The i<sup>th</sup> occurrence of a letter X in the last column corresponds to the i<sup>th</sup> occurrence of X in the first column.

# BWT Search

BWTSearch(aba)          Start from the **end** of the pattern

LF Mapping

Step 1: Find the range of "a"s in the first column

BWT(unabashable)

Step 2: Look at the same range in the last column.

Step 3: "b" is the next pattern character. Set B = the LF mapping entry for b in the first row of the range.

Set E = the LF mapping entry for b in the last + 1 row of the range.

Step 4: Find the range for "b" in the first row, and use B and E to find the right subrange within the "b" range.

| BWT | $ | a | b | e | h | l | n | s | u |
|---|---|---|---|---|---|---|---|---|---|
| $unabashable | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| abashable$un | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| able$unabash | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| ashable$unab | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| bashable$una | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| ble$unabasha | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| e$unabashabl | 0 | 2 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| hable$unabas | 0 | 2 | 1 | 1 | 1 | 1 | 1 | 0 | 0 |
| le$unabashab | 0 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 0 |
| nabashable$u | 0 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 0 |
| shable$unaba | 0 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 1 |
| unabashable$ | 0 | 3 | 2 | 1 | 1 | 1 | 1 | 1 | 1 |
|  | 1 | 3 | 2 | 1 | 1 | 1 | 1 | 1 | 1 |

# BWT Searching Example 2

pattern = "bana"

**a**

| | $ a b n |
|---|---|
| $bananna | 0 0 0 0 |
| a$banann | 0 1 0 0 |
| ananna$b | 0 1 0 1 |
| anna$ban | 0 1 1 1 |
| bananna$ | 0 1 1 2 |
| na$banan | 1 1 1 2 |
| nanna$ba | 1 1 1 3 |
| nna$bana | 1 2 1 3 |
| | 1 3 1 3 |

**n**

| | $ a b n |
|---|---|
| $bananna | 0 0 0 0 |
| a$banann | 0 1 0 0 |
| ananna$b | 0 1 0 1 |
| anna$ban | 0 1 1 1 |
| bananna$ | 0 1 1 2 |
| na$banan | 1 1 1 2 |
| nanna$ba | 1 1 1 3 |
| nna$bana | 1 2 1 3 |
| | 1 3 1 3 |

(B,E) = 0, 2

**n**

| | $ a b n |
|---|---|
| $bananna | 0 0 0 0 |
| a$banann | 0 1 0 0 |
| ananna$b | 0 1 0 1 |
| anna$ban | 0 1 1 1 |
| bananna$ | 0 1 1 2 |
| na$banan | 1 1 1 2 |
| nanna$ba | 1 1 1 3 |
| nna$bana | 1 2 1 3 |
| | 1 3 1 3 |

**a**

| | $ a b n |
|---|---|
| $bananna | 0 0 0 0 |
| a$banann | 0 1 0 0 |
| ananna$b | 0 1 0 1 |
| anna$ban | 0 1 1 1 |
| bananna$ | 0 1 1 2 |
| na$banan | 1 1 1 2 |
| nanna$ba | 1 1 1 3 |
| nna$bana | 1 2 1 3 |
| | 1 3 1 3 |

(B,E) = 1, 2

**a**

| | $ a b n |
|---|---|
| $bananna | 0 0 0 0 |
| a$banann | 0 1 0 0 |
| ananna$b | 0 1 0 1 |
| anna$ban | 0 1 1 1 |
| bananna$ | 0 1 1 2 |
| na$banan | 1 1 1 2 |
| nanna$ba | 1 1 1 3 |
| nna$bana | 1 2 1 3 |
| | 1 3 1 3 |

(B,E) = 0, 1

**b**

| | $ a b n |
|---|---|
| $bananna | 0 0 0 0 |
| a$banann | 0 1 0 0 |
| ananna$b | 0 1 0 1 |
| anna$ban | 0 1 1 1 |
| bananna$ | 0 1 1 2 |
| na$banan | 1 1 1 2 |
| nanna$ba | 1 1 1 3 |
| nna$bana | 1 2 1 3 |
| | 1 3 1 3 |

# BWT Searching Notes

- Don't have to store the LF mapping. A more complex algorithm lets you compute it in $O(1)$ time on the fly with only a little bit of storage.

- To find the range in the first column corresponding to a character:

  - Pre-compute array $C[c]$ = # of occurrences in the string of characters lexicographically < c.

  - Then start of the "a" range, for example, is: $C["a"] + 1$.

- Running time: $O(|pattern|)$

  - Finding the range in the first column takes $O(1)$ time using the C array.

  - Updating the range takes $O(1)$ time using the LF mapping.

# Relationship Between BWT and Suffix Arrays

s = appellee$
123456789

| BWT matrix | The suffixes |
|---|---|
| $appellee | $ |
| appellee$ | appellee$ |
| e$appelle | e$ |
| ee$appell | ee$ |
| ellee$app | ellee$ |
| lee$appel | lee$ |
| llee$appe | llee$ |
| pellee$ap | pellee$ |
| ppellee$a | ppellee$ |

These are still in sorted order because "$" comes before everything else

Suffix array (start position for the suffixes):
9
1
8
7
4
6
5
3
2

— subtract 1 →

s[9-1] = e
s[1-1] = $
s[8-1] = e
s[7-1] = l
s[4-1] = p
s[6-1] = l
s[5-1] = e
s[3-1] = p
s[2-1] = a

**BWT matrix**

The suffixes are obtained by deleting everything after the $

Suffix array (start position for the suffixes)

Suffix position - 1 = the position of the last character of the BWT matrix
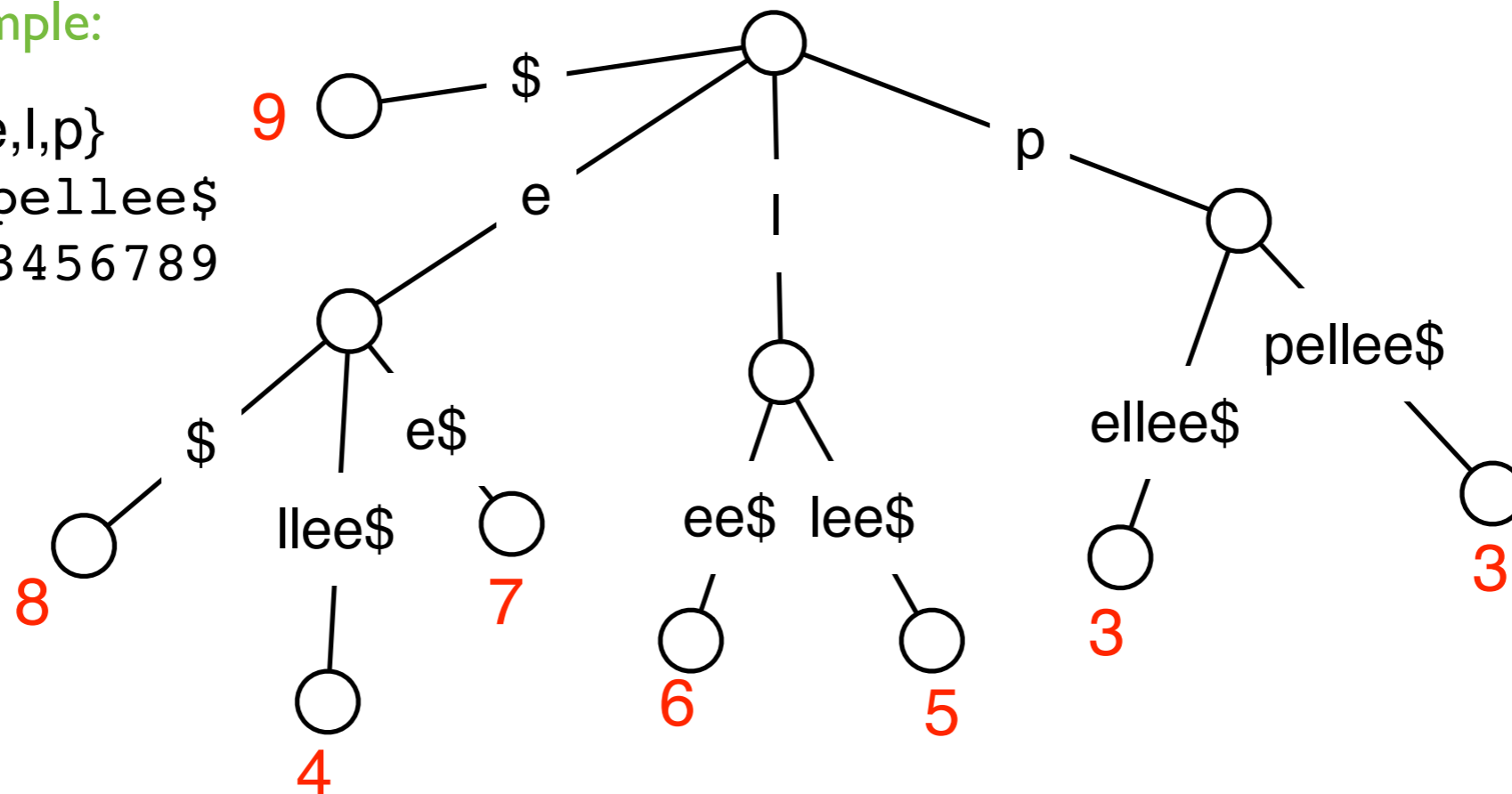
($ is a special case)

# Relationship Between BWT and Suffix Trees

- Remember: Suffix Array = suffix numbers obtained by traversing the leaf nodes of the (ordered) Suffix Tree from left to right.

- Suffix Tree ⇒ Suffix Array ⇒ BWT.

Ordered suffix tree for previous example:

$\Sigma = \{\$,e,l,p\}$

```
s = appellee$
    123456789
```

# Computing BWT in O(n) time

- Easy $O(n^2 \log n)$-time algorithm to compute the BWT (create and sort the BWT matrix explicitly).

- Several direct $O(n)$-time algorithms for BWT.
  These are space efficient.

- Also can use suffix arrays or trees:

  Compute the suffix array, use correspondence between suffix array and BWT to output the BWT.

  $O(n)$-time and $O(n)$-space, but the constants are large.

# Recap

BWT useful for searching and compression.

BWT is *invertible*: given the BWT of a string, the string can be reconstructed!

BWT is computable in O(n) time.

Close relationships between Suffix Trees, Suffix Arrays, and BWT:

- Suffix array = order of the suffix numbers of the suffix tree, traversed left to right

- BWT = letters at positions given by the suffix array entries - 1