# Burrows-Wheeler Transform and Suffix Arrays

Pavel Pevzner

Department of Computer Science and Engineering
University of California at San Diego

## Algorithms on Strings
## Data Structures and Algorithms

This slide desk is incomplete.
For the complete set of frames,
please see our videos in the
[Algorithms on Strings](#) course on [Coursera](#)
([Algorithms and Data Structures](#) Specialization)

# Outline

- **Burrows-Wheeler Transform**

- Inverting Burrows-Wheeler Transform

- Using BWT for Pattern Matching

- Suffix Arrays

- Approximate Pattern Matching

# Text Compression by Run-Length Encoding

- **Run-length encoding** compresses a run of *n* identical symbols:

*Text*

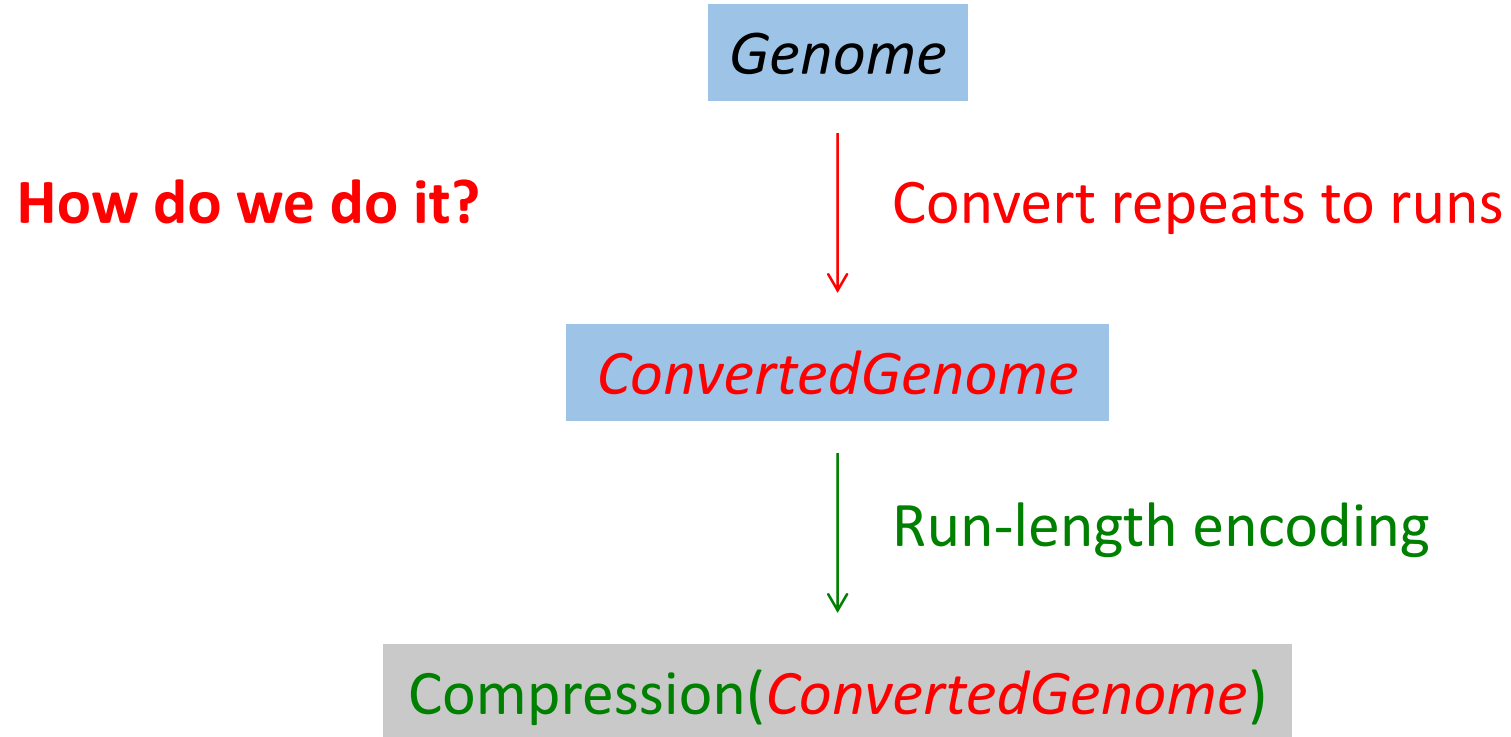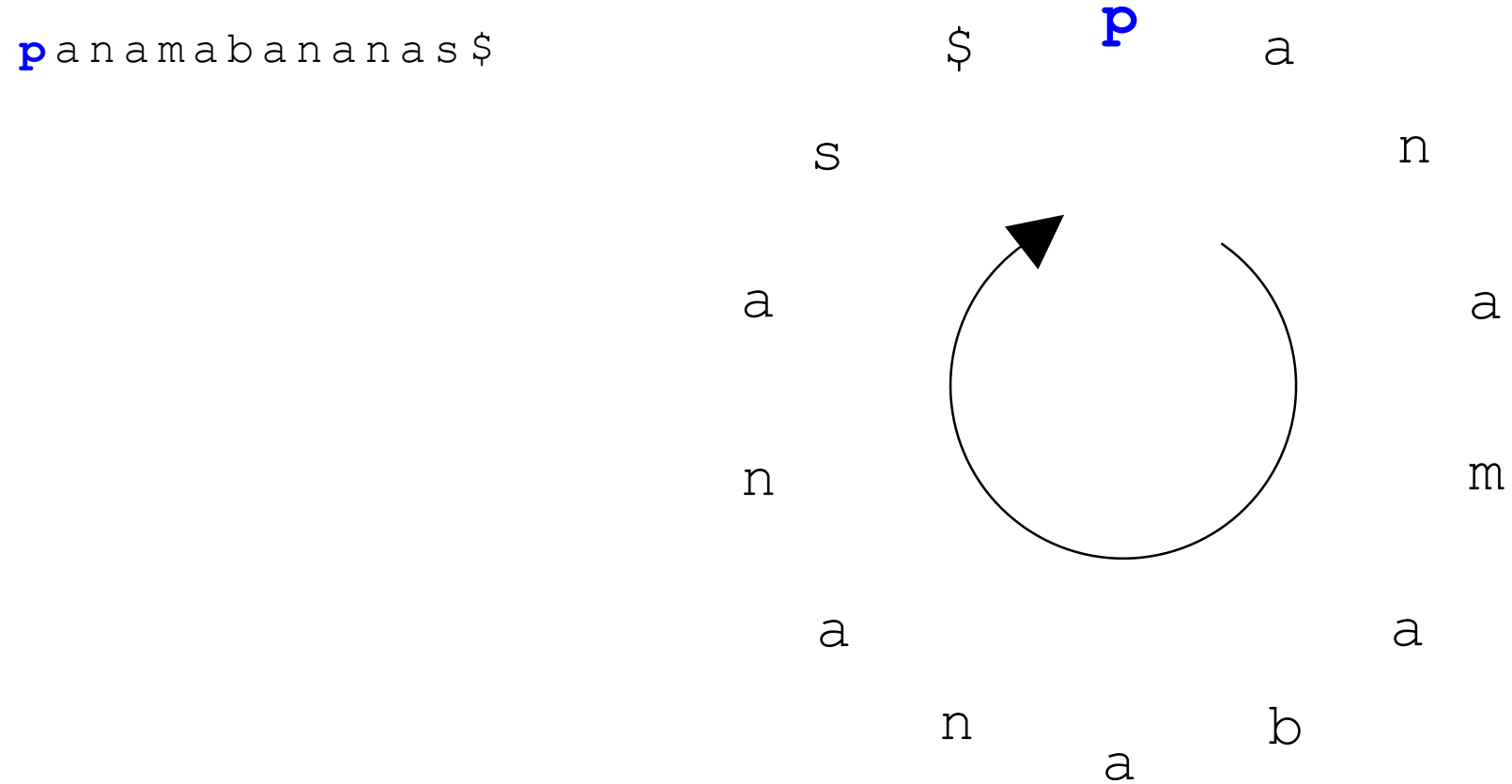GGGGGGGGGGCCCCCCCCCCCAAAAAAATTTTTTTTTTTTTTTCCCCCG

10G11C7A15T5C1G

- genomes don't have lots of runs... but they do have lots of repeats:

ACTGACCGAAACTGAGTATCCGACTGAAACTGATCAGTACTGACATTGC
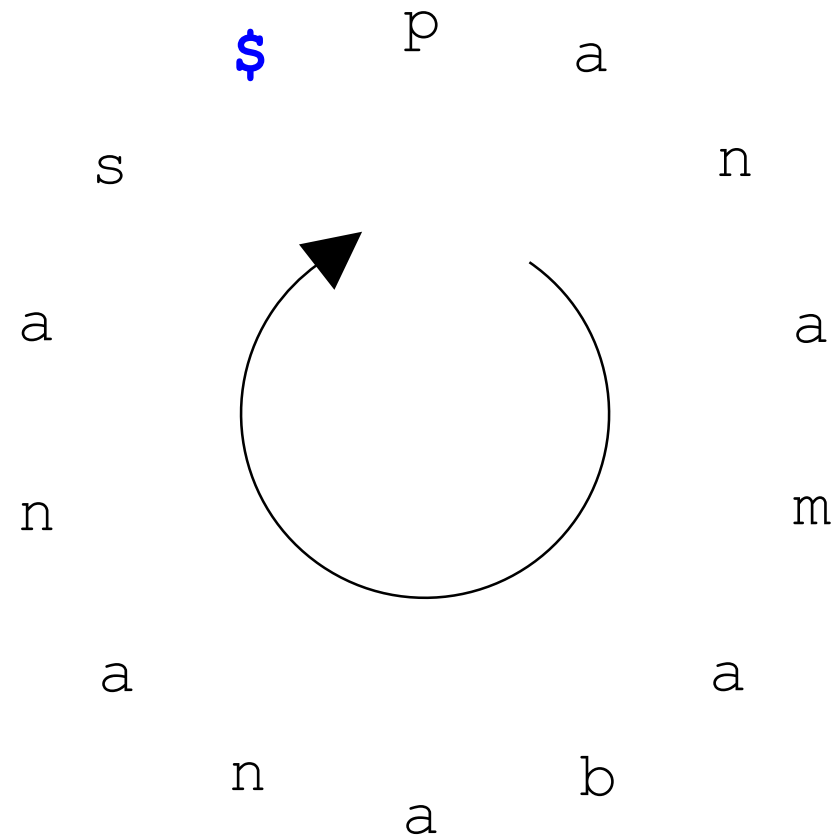
# Idea: Converting Repeats to Runs

**How do we do it?**

Genome

↓ Convert repeats to runs

*ConvertedGenome*

↓ Run-length encoding

Compression(*ConvertedGenome*)

# Forming All Cyclic Rotations of *Text*

**p**anamabananas$

$ **p** a

s n

a a

n m

a a

n b

a

# Cyclic Rotations

```
panamabananas$
$panamabananas
```

# Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
```

$ p a n a s n a m n a b a a n a

# Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
```
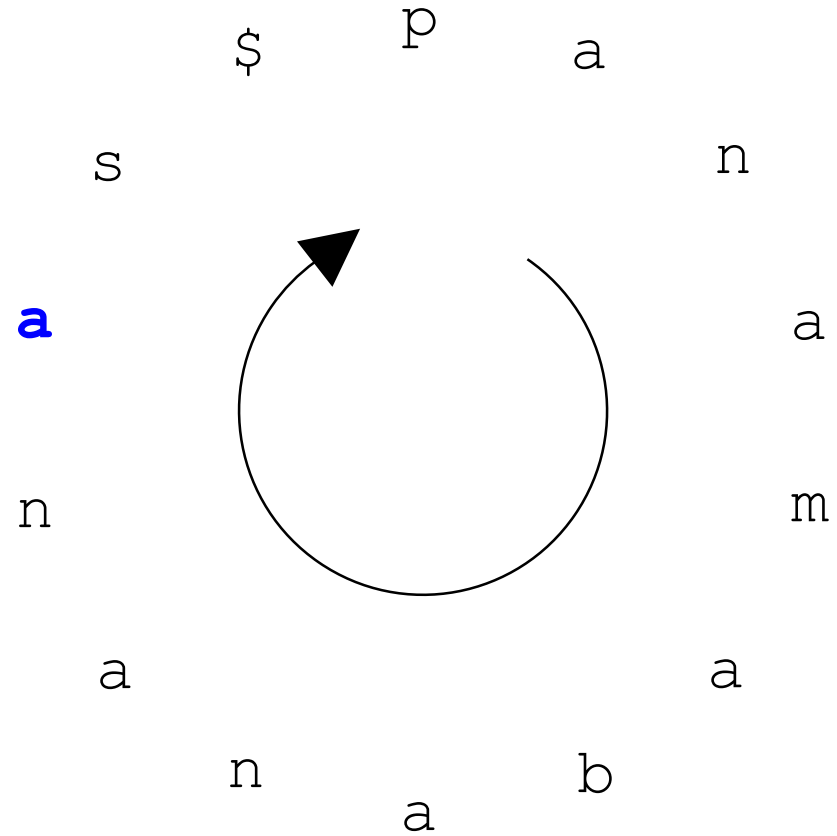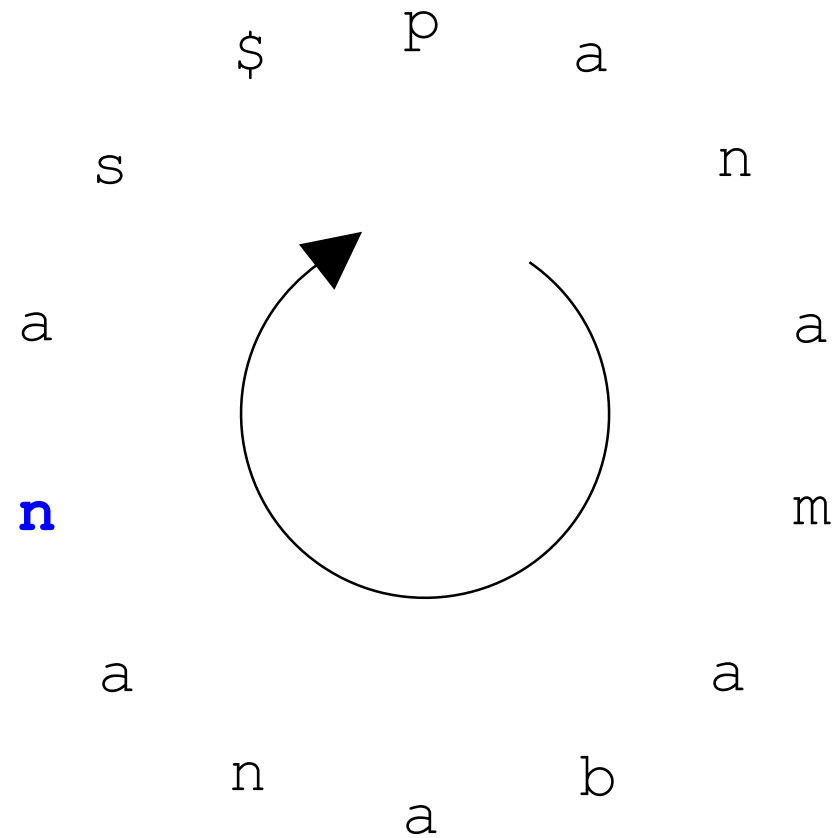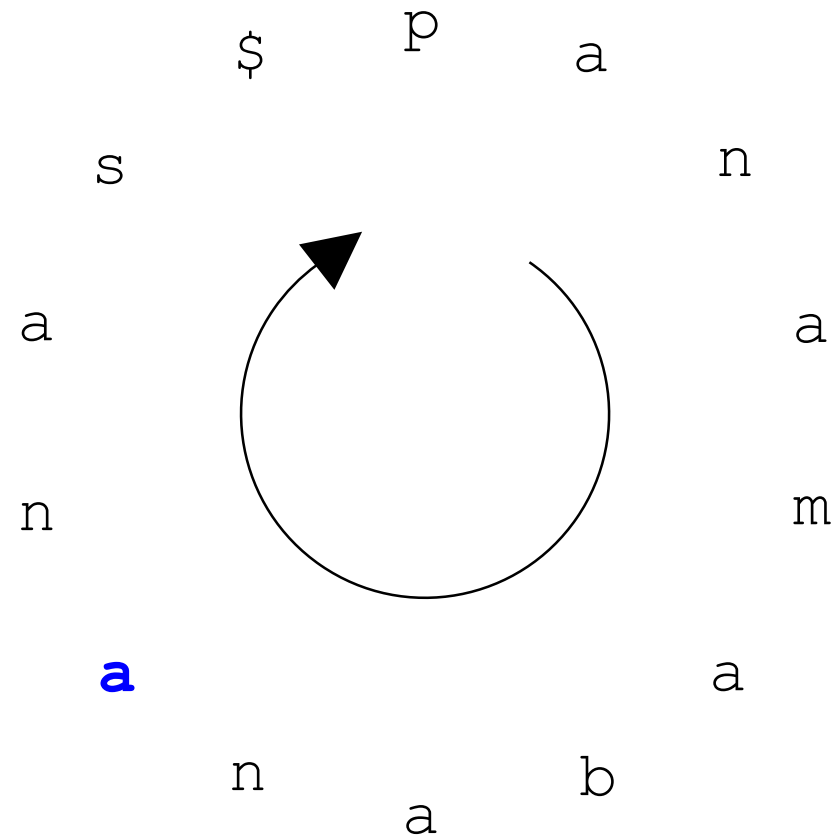
# Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
```
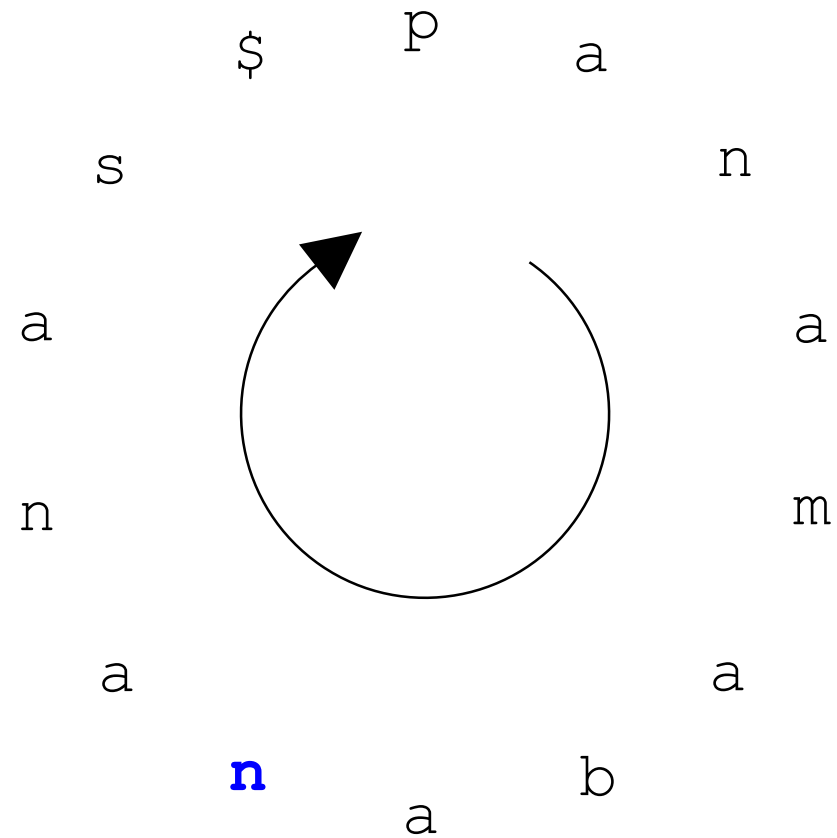
# Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
```
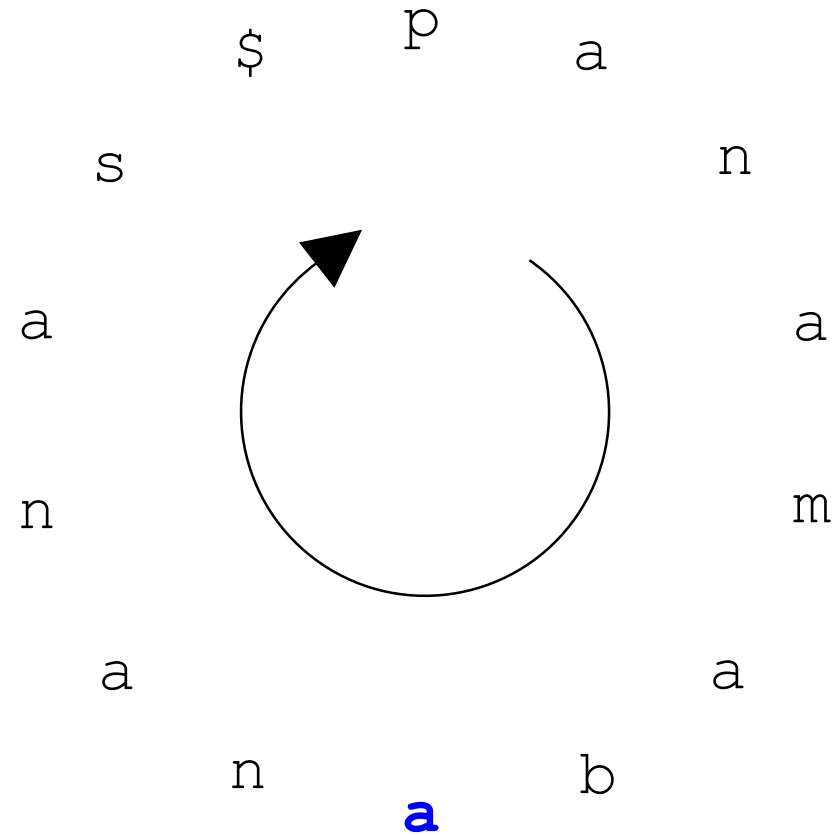
# Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
```

# Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
ananas$panamab
```

# Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
ananas$panamab
bananas$panama
```

# Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
ananas$panamab
bananas$panama
abananas$panam
```
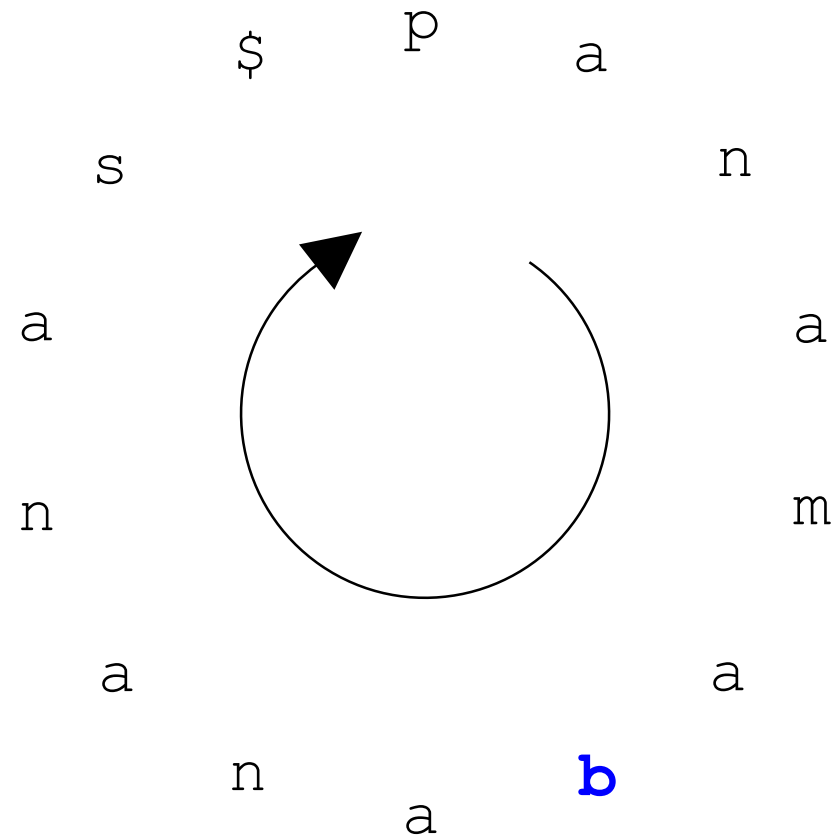
# Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
ananas$panamab
bananas$panama
abananas$panam
mabananas$pana
```

# Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
ananas$panamab
bananas$panama
abananas$panam
mabananas$pana
amabananas$pan
```
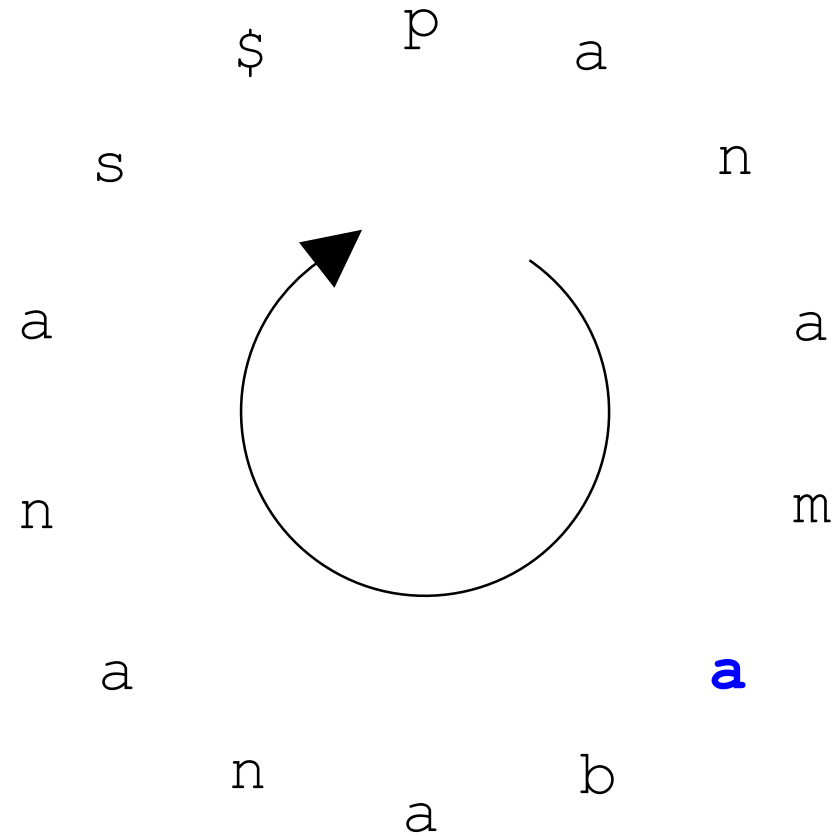
# Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
ananas$panamab
bananas$panama
abananas$panam
mabananas$pana
amabananas$pan
namabananas$pa
```
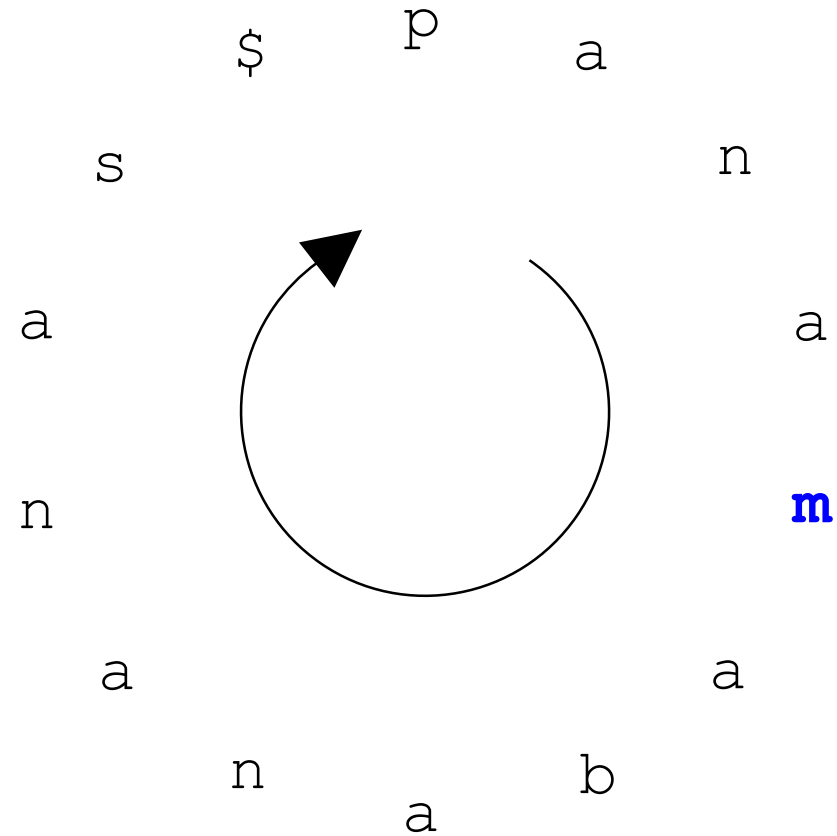
# Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
ananas$panamab
bananas$panama
abananas$panam
mabananas$pana
amabananas$pan
namabananas$pa
anamabananas$p
```

# Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
ananas$panamab
bananas$panama
abananas$panam
mabananas$pana
amabananas$pan
namabananas$pa
anamabananas$p
```
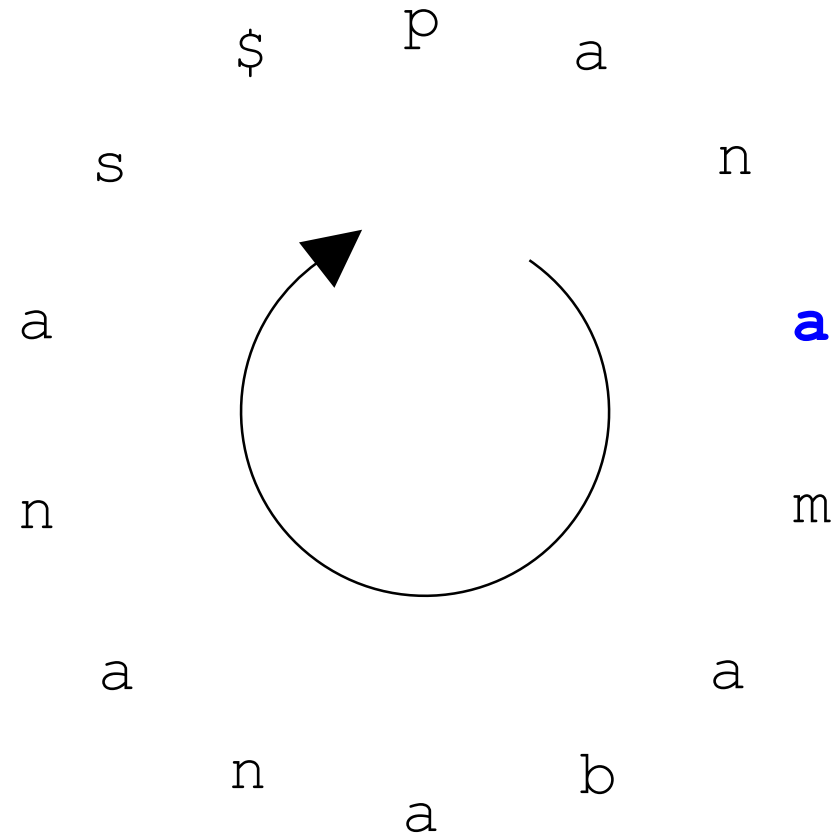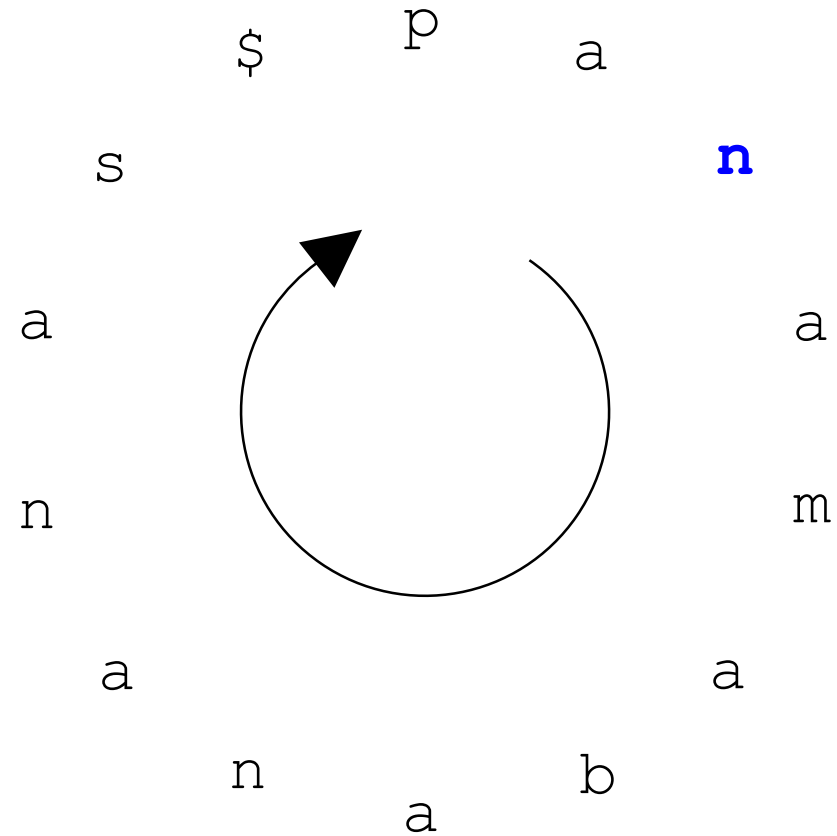
$
p
a
n
a
m
a
n
a
b
a
n
a
s

# Sorting Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
ananas$panamab
bananas$panama
abananas$panam
mabananas$pana
amabananas$pan
namabananas$pa
anamabananas$p
```

**$**panamabananas
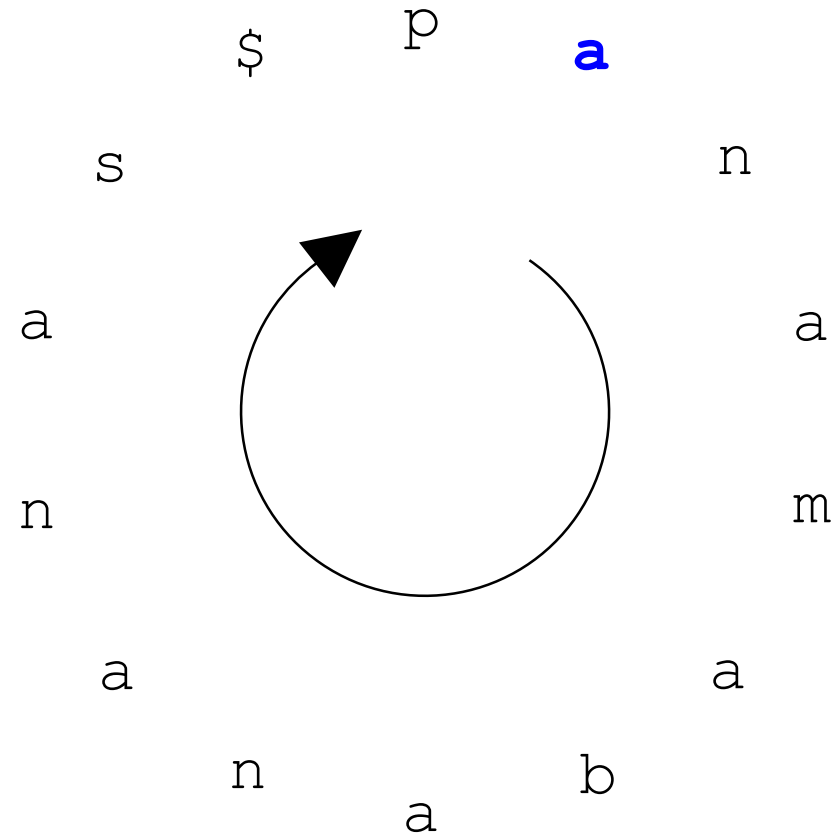
Sort the strings
lexicographically
($ comes first)

# Sorting Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
ananas$panamab
bananas$panama
abananas$panam
mabananas$pana
amabananas$pan
namabananas$pa
anamabananas$p
```
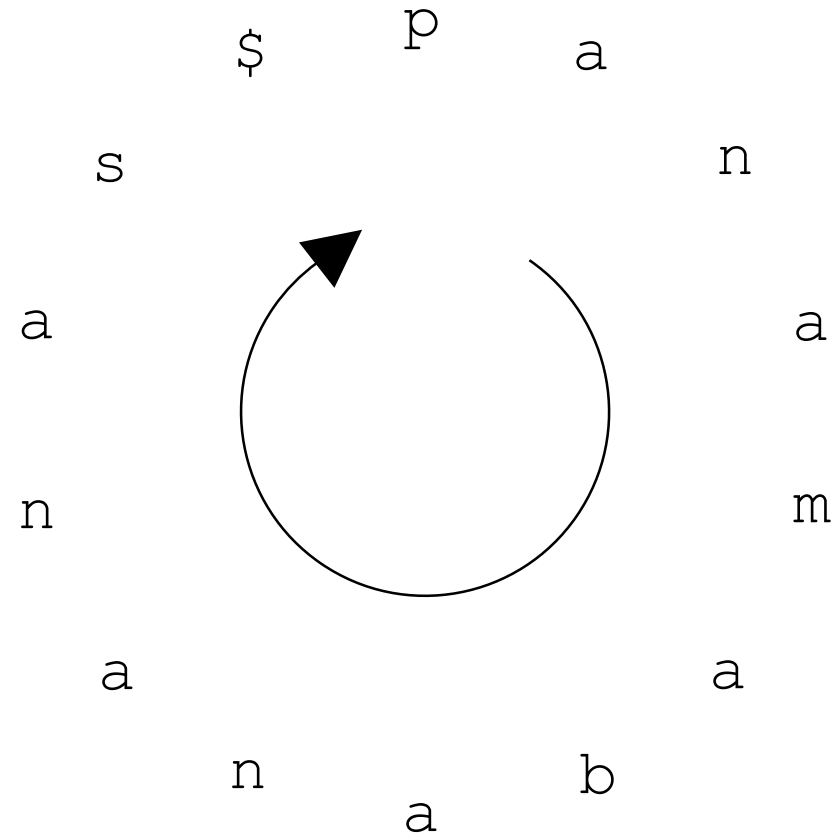
**$**panamabananas
**a**bananas$panam

# Sorting Cyclic Rotations

panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
ananas$panamab
bananas$panama
abananas$panam
mabananas$pana
amabananas$pan
namabananas$pa
anamabananas$p

**$**panamabananas
**a**bananas$panam
**am**abananas$pan

Sort the strings
lexicographically
($ comes first)

# Sorting Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
ananas$panamab
bananas$panama
abananas$panam
mabananas$pana
amabananas$pan
namabananas$pa
anamabananas$p
```

**$**panamabananas
**a**bananas$panam
**am**abananas$pan
**anam**abananas$p

Sort the strings
lexicographically
($ comes first)

# Sorting Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
ananas$panamab
bananas$panama
abananas$panam
mabananas$pana
amabananas$pan
namabananas$pa
anamabananas$p
```

**$**panamabananas
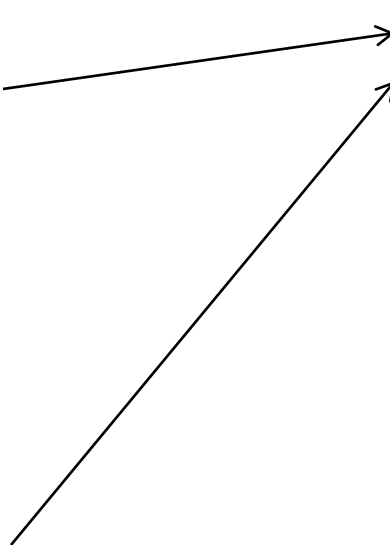**a**bananas$panam
**am**abananas$pan
**anam**abananas$p
**anan**as$panamab

Sort the strings
lexicographically
($ comes first)

# Sorting Cyclic Rotations

```
panamabananas$
$panamabananas
s$panamabanana
as$panamabanan
nas$panamabana
anas$panamaban
nanas$panamaba
ananas$panamab
bananas$panama
abananas$panam
mabananas$pana
amabananas$pan
namabananas$pa
anamabananas$p
```
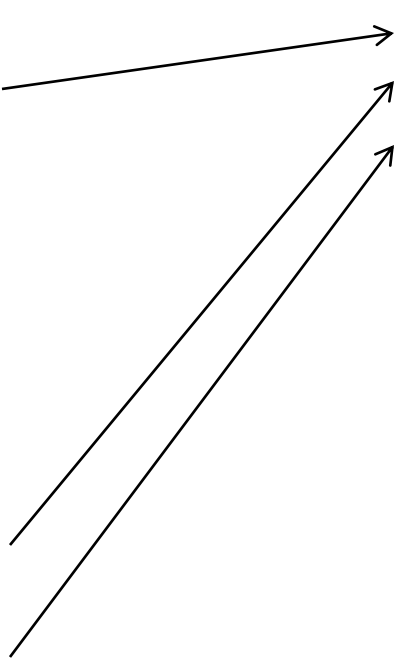
→

```
$panamabananas
abananas$panam
amabananas$pan
anamabananas$p
ananas$panamab
anas$panamaban
as$panamabanan
bananas$panama
mabananas$pana
namabananas$pa
nanas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```

Sort the strings
lexicographically
($ comes first)

# BWT(panamabananas$)=smnpbnnaaaaa$a

```
panamabananas$            $panamabanana s
$panamabananas            abananas$panam m
s$panamabanana            amabananas$pan n
as$panamabanan            anamabananas$p p
nas$panamabana            ananas$panamab b
anas$panamaban            anas$panamaban n
nanas$panamaba    ———→    as$panamabanan n
ananas$panamab            bananas$panama a
bananas$panama            mabananas$pana a
abananas$panam            namabananas$pa a
mabananas$pana            nanas$panamab a
amabananas$pan            nas$panamaban a
namabananas$pa            panamabananas $
anamabananas$p            s$panamabanan a
```

All cyclic rotations of                 **Burrows-Wheeler Transform (BWT)**:
"panamabananas$"                         Last column = **smnpbnnaaaaa$a**

# BWT(panamabananas$)=smnpbnnaaaaa$a

```
panamabananas$              $panamabananas
$panamabananas              abananas$panam
s$panamabanana              amabananas$pan
as$panamabanan              anamabananas$p
nas$panamabana              ananas$panamab
anas$panamaban              anas$panamaban
nanas$panamaba   ------>    as$panamabanan
ananas$panamab              bananas$panama
bananas$panama              mabananas$pana
abananas$panam              namabananas$pa
mabananas$pana              nanas$panamaba
amabananas$pan              nas$panamabana
namabananas$pa              panamabananas$
anamabananas$p              s$panamabanana
```
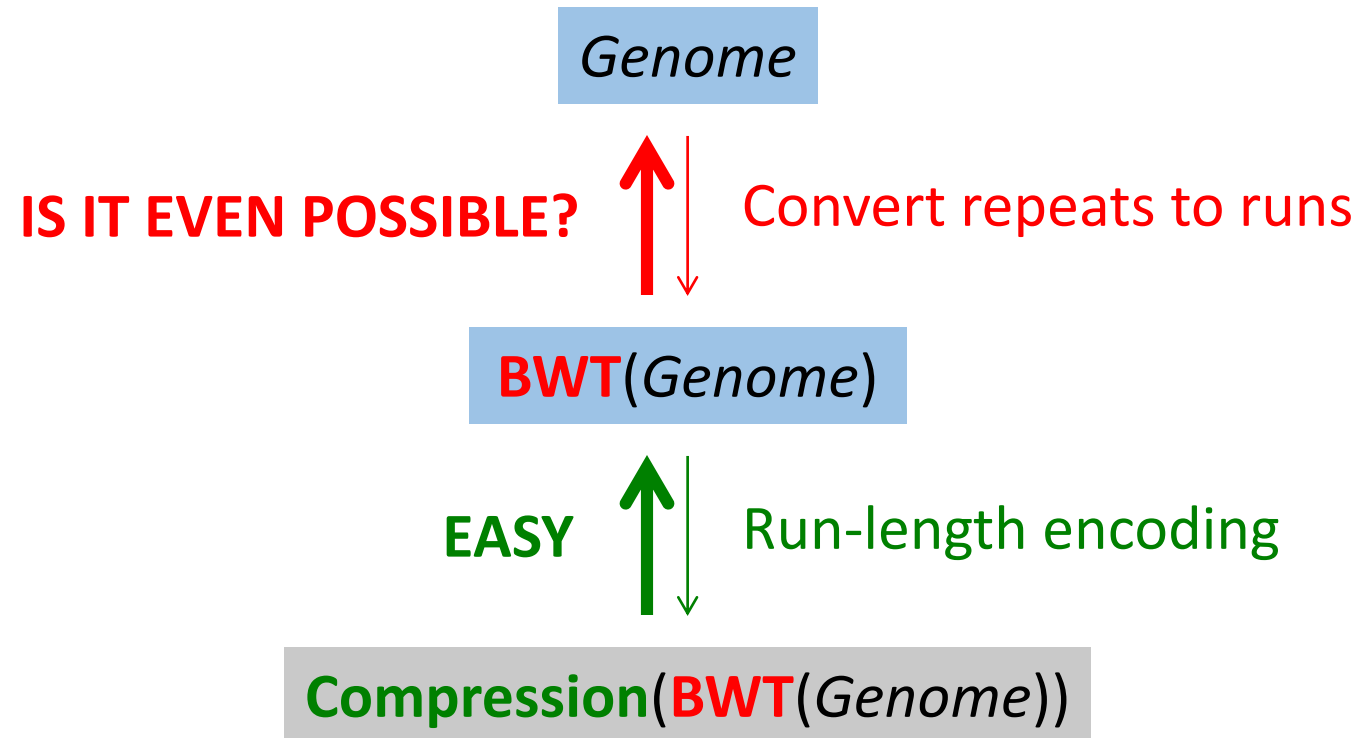
All cyclic rotations of
"panamabananas$"

**Burrows-Wheeler Transform (BWT)**:
Last column = **smnpbnnaaaaa$a**

# Applying BWT to the Double Helix Paper by Watson&Crick

nd Corey (1).  They kindly made their manuscript availa ....... **a**
nd criticism, especially on interatomic distances.  We  ...... **a**
nd cytosine.  The sequence of bases on a single chain d ...... **a**
nd experimentally (3,4) that the ratio of the amounts o ...... u
nd for this reason we shall not comment on it.  We wish ...... **a**
nd guanine (purine) with cytosine (pyrimidine).  In oth ...... **a**
nd ideas of Dr.  M. H. F. Wilkins, Dr.  R. E. Franklin  ...... **a**
nd its water content is rather high.  At lower water co ...... **a**
nd pyrimidine bases.  The planes of the bases are perpe ...... **a**
nd stereochemical arguments.  It has not escaped our no ...... **a**
nd that only specific pairs of bases can bond together  ...... u
nd the atoms near it is close to Furberg's 'standard co ...... **a**
nd the bases on the inside, linked together by hydrogen ...... **a**
nd the bases on the outside.  In our opinion, this stru ...... **a**
nd the other a pyrimidine for bonding to occur.  The hy ...... **a**
nd the phosphates on the outside.  The configuration of ...... **a**
nd the ration of guanine to cytosine, are always very c ...... **a**
nd the same axis (see diagram).  We have made the usual ...... u
nd their co-workers at King's College, London.  One of  ...... **a**

"and" is a frequent repeat in English texts

# Going Back From BWT(*Genome*) to *Genome*

*Genome*

**IS IT EVEN POSSIBLE?**  Convert repeats to runs

**BWT**(*Genome*)

**EASY**  Run-length encoding

**Compression**(**BWT**(*Genome*))

# Outline

- Burrows-Wheeler Transform

- **Inverting Burrows-Wheeler Transform**

- Using BWT for Pattern Matching

- Suffix Arrays

- Approximate Pattern Matching

# Reconstructing `banana` **from** `annb$aa`

```
$banana
a$banan
ana$ban
anana$b
banana$
na$bana
nana$ba
```

# Reconstructing banana

**$**banan**a**
**a**$bana**n**
**a**na$ba**n**
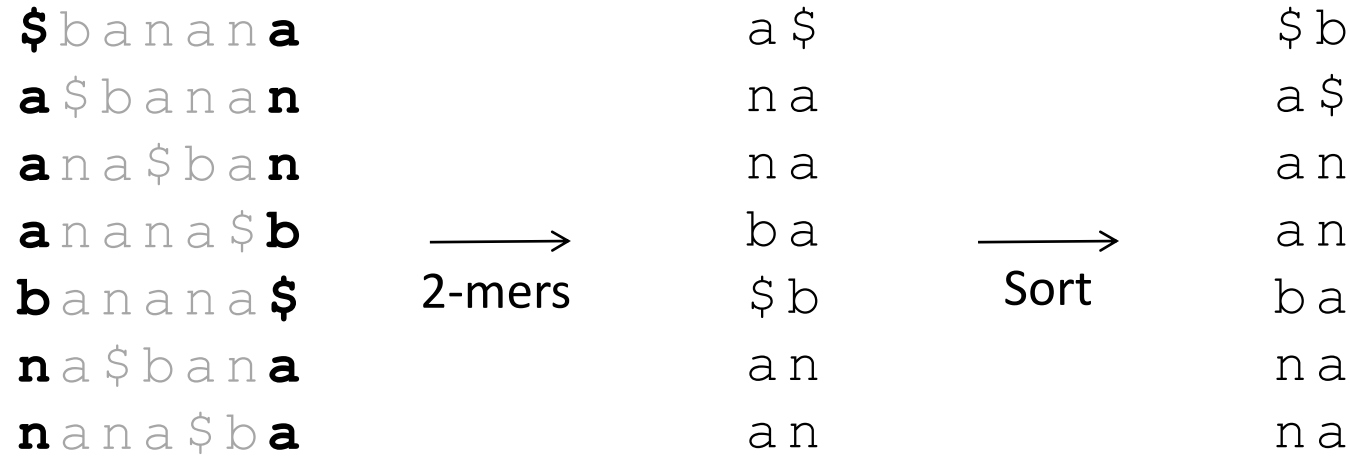**a**nana$**b**
**b**anana**$**
**n**a$bana**a**
**n**ana$b**a**

- Sorting all elements of "annb$aa" gives first column of BWT matrix.

# Reconstructing `banana`

```
$banana          a$
a$banan          na
ana$ban          na
anana$b    ⟶     ba
banana$   2-mers  $b
na$bana          an
nana$ba          an
```

- We now know 2-mer composition of the circular string `banana$`

# Reconstructing `banana`

```
$banana        a$              $b
a$banan        na              a$
ana$ban        na              an
anana$b    →   ba          →   an
banana$   2-mers  $b    Sort    ba
na$bana        an              na
nana$ba        an              na
```

- We now know 2-mer composition of the circular string `banana$`

- Sorting gives us the first 2 columns of the matrix.

# Reconstructing banana

```
$banana        a$              $b
a$banan        na              a$
ana$ban        na              an
anana$b  ──→   ba    ──→       an
banana$ 2-mers $b   Sort       ba
na$bana        an              na
nana$ba        an              na
```

- We now know 2-mer composition of the circular string `banana$`
- Sorting gives us the first 2 columns of the matrix.

# Reconstructing banana

```
$banana
a$banan
ana$ban
anana$b
banana$
na$bana
nana$ba
```

# Reconstructing banana

```
$banana          a$b
a$bann           na$
ana$ban          nan
anana$b          ban
banana$    ⟶     $ba
na$bana    3-mers ana
nana$ba          ana
```

- We now know 3-mer composition of the circular string `banana$`

# Reconstructing banana

```
$banana            a$b              $ba
a$banan            na$              a$b
ana$ban            nan              ana
anana$b    ──────> ban    ──────>   ana
banana$    3-mers  $ba    Sort      ban
na$bana            ana              na$
nana$ba            ana              nan
```

- We now know 3-mer composition of the circular string `banana$`

- Sorting gives us the first 3 columns of the matrix.

# Reconstructing banana

```
$banana          a$b              $ba
a$banan          na$              a$b
ana$ban          nan              ana
anana$b   3-mers  ban    Sort     ana
banana$          $ba              ban
na$bana          ana              na$
nana$ba          ana              nan
```

- We now know 3-mer composition of the circular string `banana$`
- Sorting gives us the first 3 columns of the matrix.

# Reconstructing banana

```
$ba nan a
a$b ana n
ana $ba n
ana na$ b
ban ana $
na$ ban a
nan a$b a
```

# Reconstructing banana

```
$banana          a$ba
a$banan          na$b
ana$ban          nana
anana$b          bana
banana$          $ban
na$bana          ana$
nana$ba          anan
```

→ 4-mers

- We now know 4-mer composition of the circular string banana$

# Reconstructing `banana`

```
$banana        a$ba           $ban
a$banan        na$b           a$ba
ana$ban        nana           ana$
anana$b  ──────→ bana  ──────→ anan
banana$  4-mers  $ban   Sort   bana
na$bana        ana$           na$b
nana$ba        anan           nana
```

- We now know 4-mer composition of the circular string `banana$`
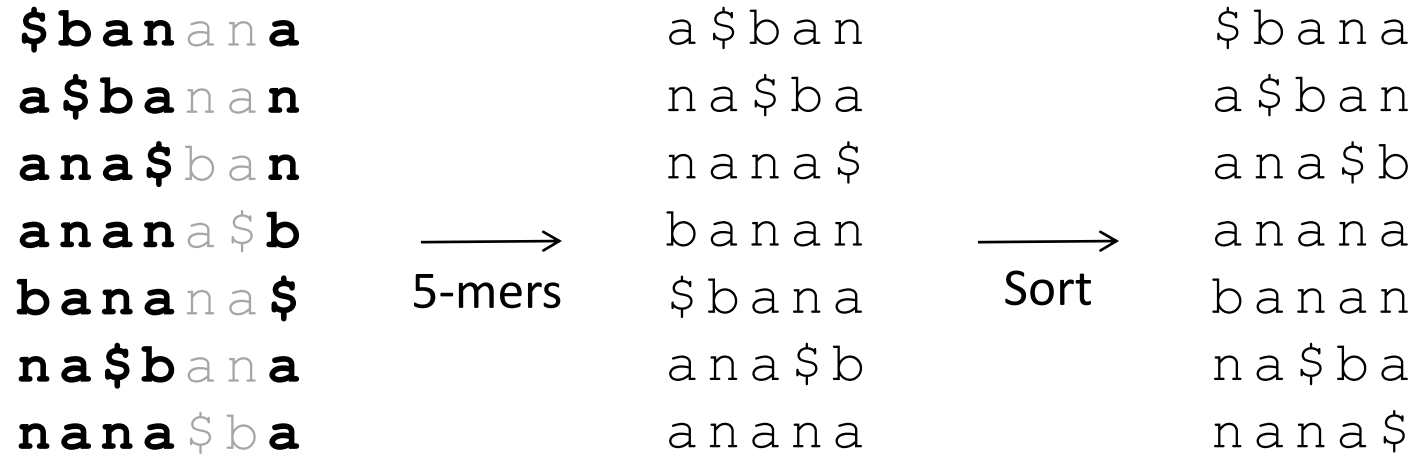- Sorting gives us the first 4 columns of the matrix.

# Reconstructing `banana`

```
$banana        a$ba              $ban
a$banan        na$b              a$ba
ana$ban        nana              ana$
anana$b    →   bana     →        anan
banana$      4-mers   $ban   Sort  bana
na$bana        ana$              na$b
nana$ba        anan              nana
```

- We now know 4-mer composition of the circular string `banana$`

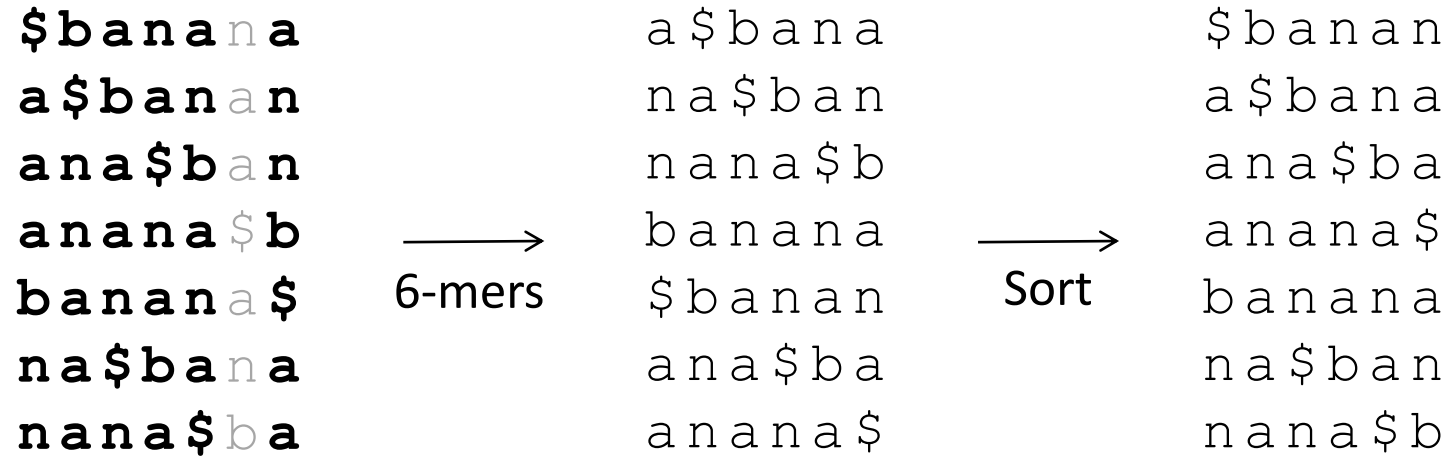- Sorting gives us the first 4 columns of the matrix.

# Reconstructing banana

```
$ban ana
a$ba nan
ana$ban
anana$b
bana na$
na$bana
nana$ba
```

# Reconstructing banana

$banana          a$ban
a$banan          na$ba
ana$ban          nana$
anana$b   ⟶     banan
banana$   5-mers  $bana
na$bana          ana$b
nana$ba          anana

- We now know 5-mer composition of the circular string banana$

# Reconstructing `banana`

```
$banana        a$ban          $bana
a$banan        na$ba          a$ban
ana$ban        nana$          ana$b
anana$b   →    banan    →     anana
banana$  5-mers $bana   Sort  banan
na$bana        ana$b          na$ba
nana$ba        anana          nana$
```

- We now know 5-mer composition of the circular string `banana$`

- Sorting gives us the first 5 columns of the matrix.

# Reconstructing banana

```
$banana        a$ban            $bana
a$banan        na$ba            a$ban
ana$ban        nana$            ana$b
anana$b   →    banan     →      anana
banana$  5-mers $bana   Sort    banan
na$bana        ana$b            na$ba
nana$ba        anana            nana$
```

- We now know 5-mer composition of the circular string `banana$`
- Sorting gives us the first 5 columns of the matrix.

# Reconstructing banana

```
$banana
a$banan
ana$ban
anana$b
banana$
na$bana
nana$ba
```

# Reconstructing banana

$banana          a$bana

a$banan          na$ban

ana$ban          nana$b

anana$b    ⟶     banana

banana$   6-mers  $banan

na$bana          ana$ba

nana$ba          anana$

- We now know 6-mer composition of the circular string banana$

# Reconstructing banana

```
$banana          a$bana          $banan
a$banan          na$ban          a$bana
ana$ban          nana$b          ana$ba
anana$b   ──→    banana    ──→   anana$
banana$   6-mers $banan    Sort  banana
na$bana          ana$ba          na$ban
nana$ba          anana$          nana$b
```

- We now know 6-mer composition of the circular string banana$
- Sorting gives us the first 6 columns of the matrix.

# Reconstructing banana

```
$banana          a$bana          $banan
a$banan          na$ban          a$bana
ana$ban          nana$b          ana$ba
anana$b   ──────▶  banana   ──────▶  anana$
banana$   6-mers  $banan   Sort    banana
na$bana          ana$ba          na$ban
nana$ba          anana$          nana$b
```

- We now know 6-mer composition of the circular string `banana$`

- Sorting gives us the first 6 columns of the matrix.

# Reconstructing `banana`

```
$banana
a$banan
ana$ban
anana$b
banana$
na$bana
nana$ba
```

- We now know the entire matrix!

# Reconstructing `banana`

`$banana`
`a$banan`
`ana$ban`
`anana$b`
`banana$`
`na$bana`
`nana$ba`

- We now know the entire matrix!

- Symbols in the first row (after $) spell **banana**.

# More Memory Issues

- Reconstructing *Text* from *BWT*(*Text*) required us to store |*Text*| cyclic rotations of |*Text*|.

```
$banana
a$banan
ana$ban
anana$b
banana$
na$bana
nana$ba
```

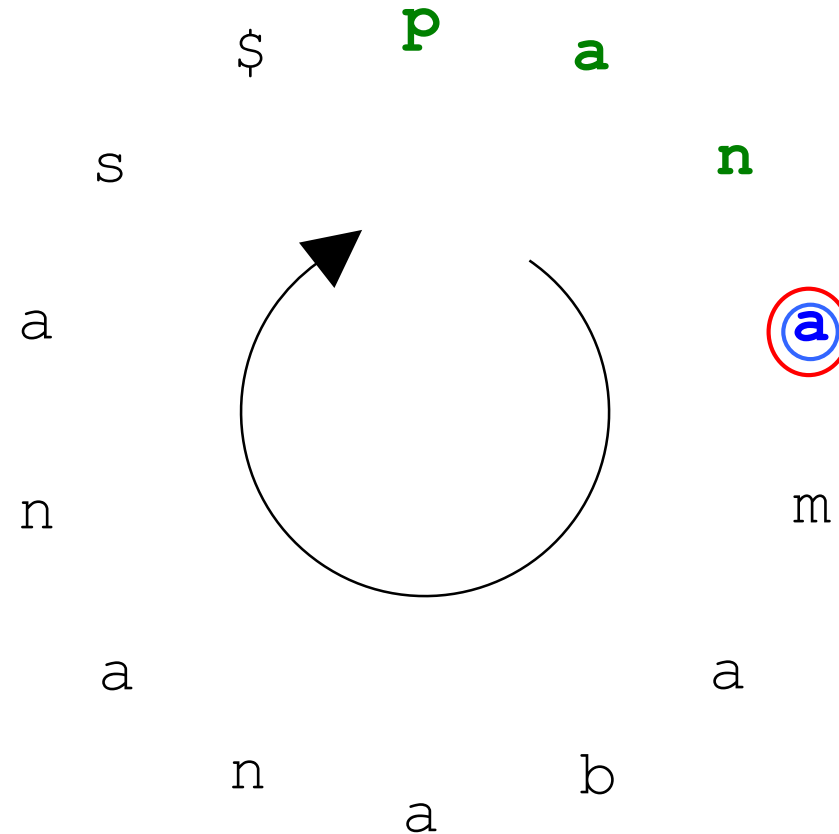- Can we invert BWT(*Text*) with less space and without |*Text*| rounds of sorting?

# A Strange Observation

```
$panamabananas
abananas$panam
amabananas$pan
anamabananas$p
ananas$panamab
anas$panamaban
as$panamabanan
bananas$panama
mabananas$pana
namabananas$pa
nanas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```
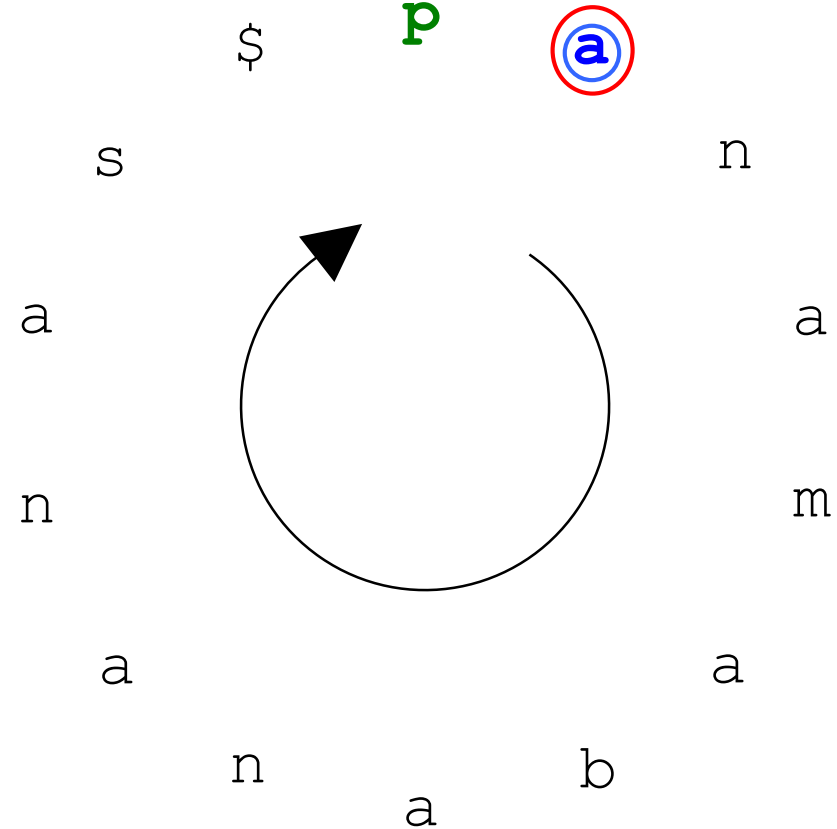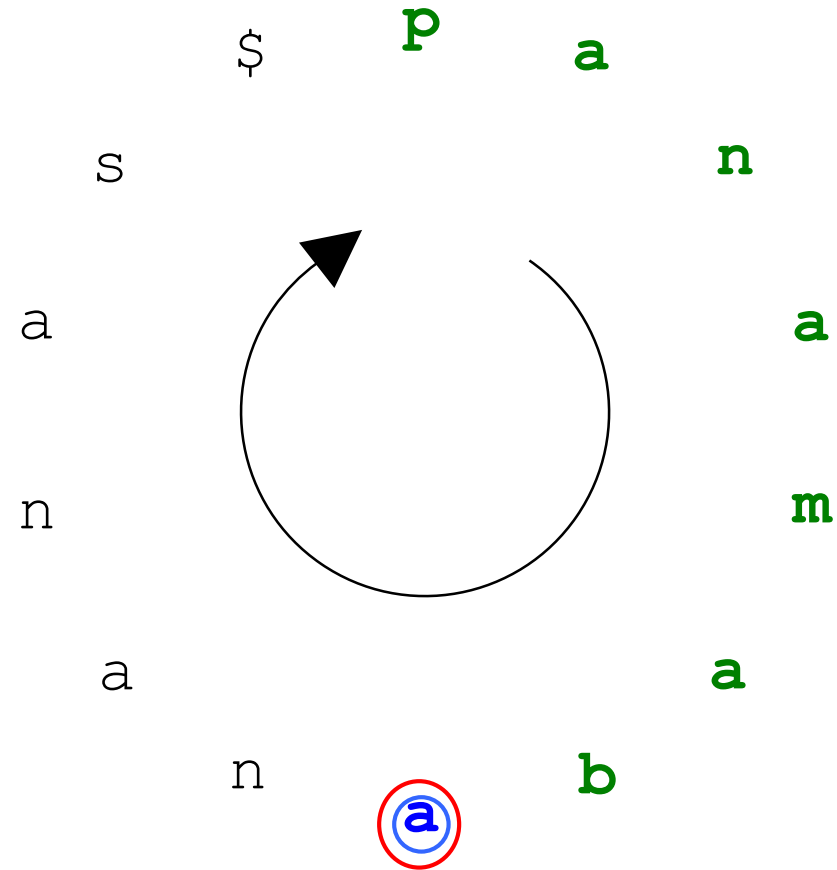
# A Strange Observation

# A Strange Observation

```
$ p a n a m a b a n a n a s
a b a n a n a s $ p a n a m
a m a b a n a n a s $ p a n
a n a m a b a n a n a s $ p
a n a n a s $ p a n a m a b
a n a s $ p a n a m a b a n
a s $ p a n a m a b a n a n
b a n a n a s $ p a n a m a
m a b a n a n a s $ p a n a
n a m a b a n a n a s $ p a
n a n a s $ p a n a m a b a
n a s $ p a n a m a b a n a
p a n a m a b a n a n a s $
s $ p a n a m a b a n a n a
```

# A Strange Observation

Where is first "a" hiding inside the circle?

```
$ p a n a m a b a n a n a s
a b a n a n a s $ p a n a m
a m a b a n a n a s $ p a n
a n a m a b a n a n a s $ p
a n a n a s $ p a n a m a b
a n a s $ p a n a m a b a n
a s $ p a n a m a b a n a n
b a n a n a s $ p a n a m a
m a b a n a n a s $ p a n a
n a m a b a n a n a s $ p a
n a n a s $ p a n a m a b a
n a s $ p a n a m a b a n a
p a n a m a b a n a n a s $
s $ p a n a m a b a n a n a
```

Where is first "a" hiding inside the circle?

# They Are Hiding at the Same Position!



1st **a** in *FirstColumn* and 1st **a** in *LastColumn*
are hiding at the same position along the cycle!

# Another Strange Observation

```
$panamabananas
abananas$panam
amabananas$pan
anamabananas$p
ananas$panamab
anas$panamaban
as$panamabanan
bananas$panama
mabananas$pana
namabananas$pa
nanas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```
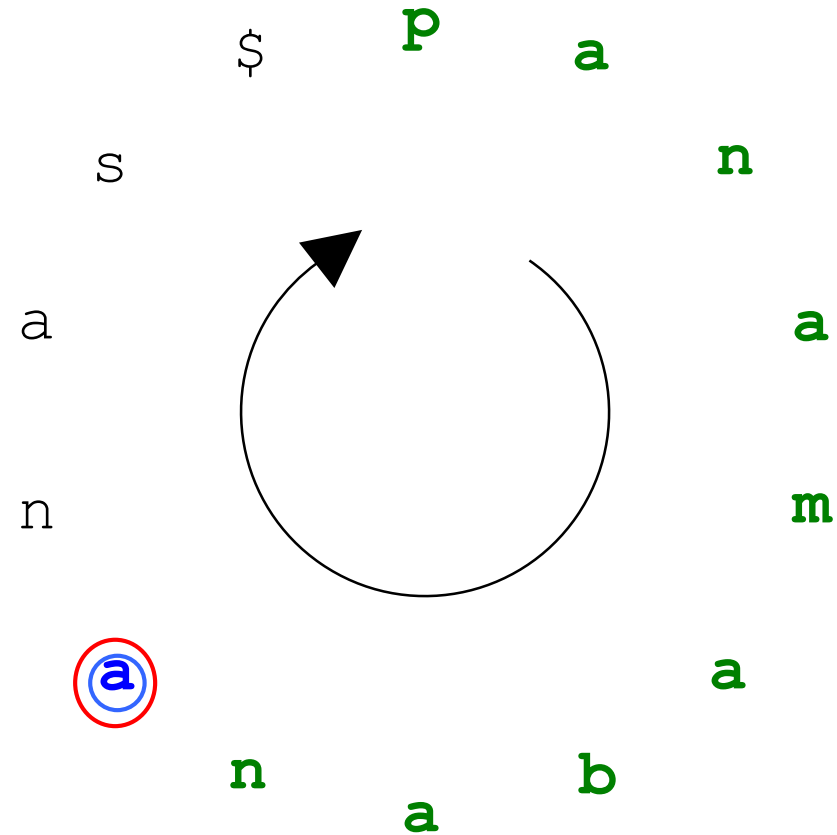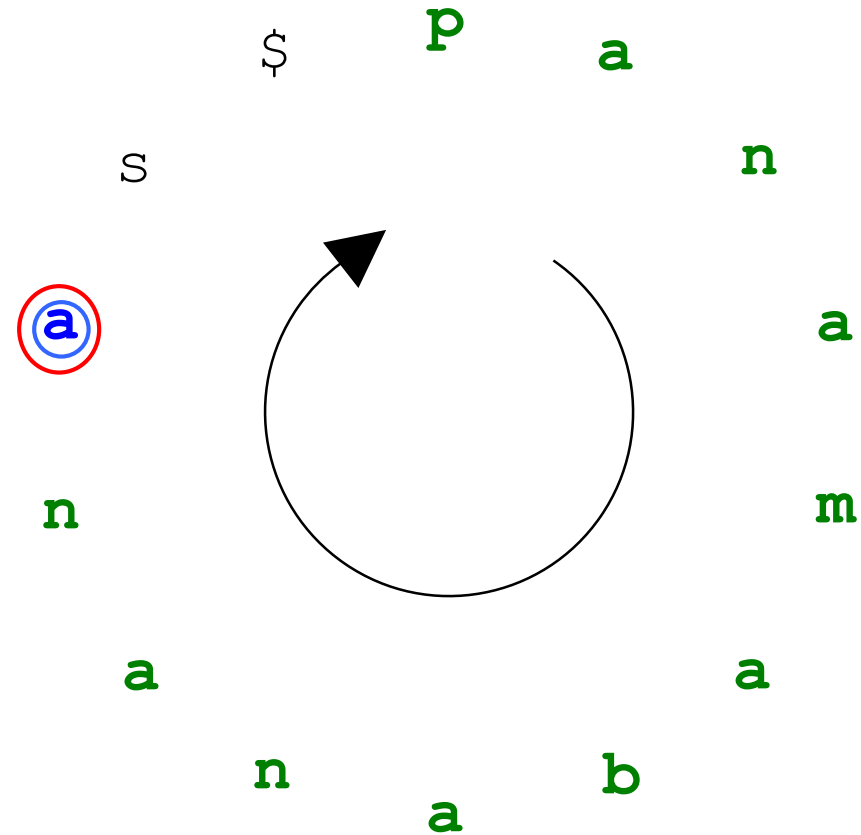


2nd **a** in *FirstColumn* and 2nd **a** in *LastColumn*
are hiding at the same position along the cycle!

# Another Strange Observation

```
$panamabananas
abananas$panam
amabananas$pan
anamabananas$p
ananas$panamab
anas$panamaban
as$panamabanan
bananas$panama
mabananas$pana
namabananas$pa
nanas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```

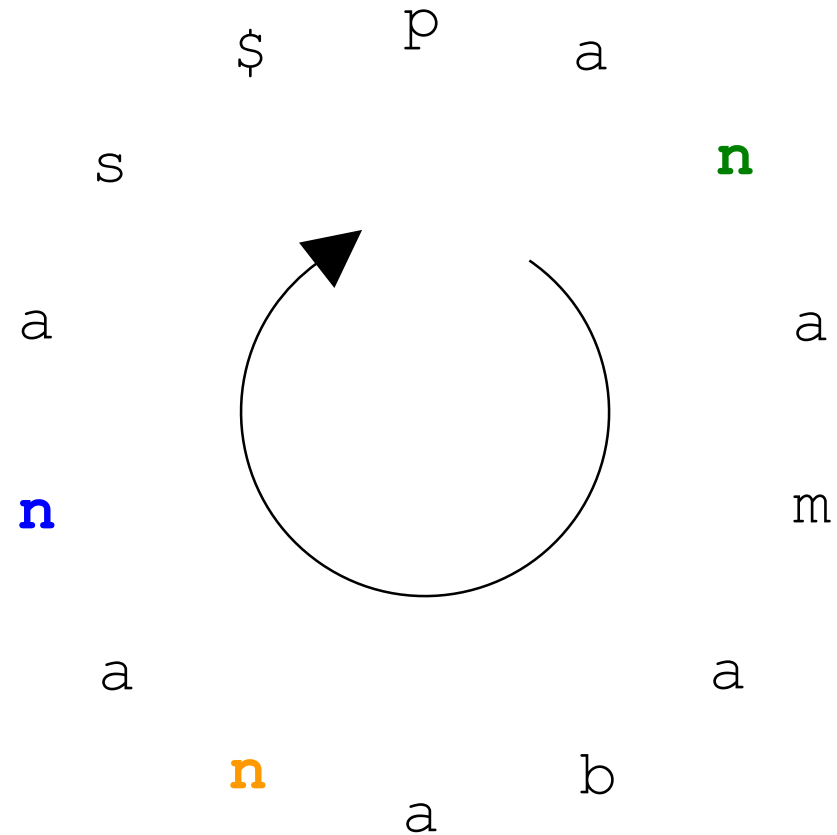# Another Strange Observation

```
$panamananas
abananas$panam
amabananas$pan
anamabananas$p
ananas$panamab
anas$panamaban
as$panamabanan
bananas$panama
mabananas$pana
namabananas$pa
nanas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```

# Another Strange Observation

# Another Strange Observation



```
$panamabananas
abananas$panam
amabananas$pan
anamabananas$p
ananas$panamab
anas$panamaban
as$panamabanan
bananas$panama
mabananas$pana
namabananas$pa
nanas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```

# Another Strange Observation

# Is It True in General?

```
  $panamabananas
1 abananas$panam
2 amabananas$pan
3 anamabananas$p
4 ananas$panamab
5 anas$panamaban
6 as$panamabanan
  bananas$panama
  mabananas$pana
  namabananas$pa
  nanas$panamaba
  nas$panamabana
  panamabananas$
  s$panamabanana
```

These strings are sorted

# Is It True in General?

```
   $panamabananas
1  abananas$panam
2  amabananas$pan
3  anamabananas$p
4  ananas$panamab
5  anas$panamaban
6  as$panamabanan
   bananas$panama
   mabananas$pana
   namabananas$pa
   nanas$panamaba
   nas$panamabana
   panamabananas$
   s$panamabanana
```

These strings are sorted

Chop off **a**

```
bananas$panam
mabananas$pan
namabananas$p
nanas$panamab
nas$panamaban
s$panamabanan
```
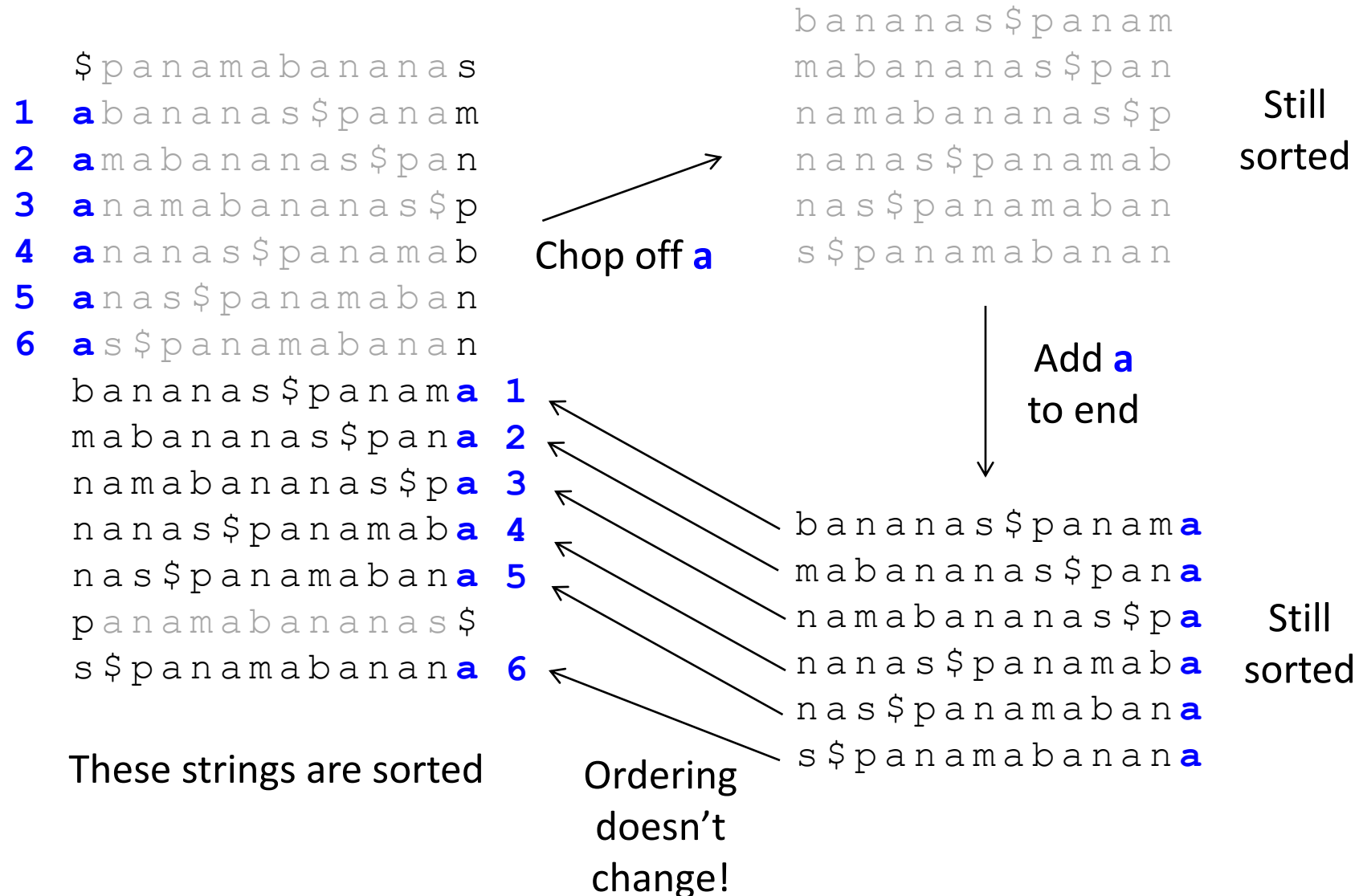
# Is It True in General?

```
  $ p a n a m a b a n a n a s
1 a b a n a n a s $ p a n a m
2 a m a b a n a n a s $ p a n
3 a n a m a b a n a n a s $ p
4 a n a n a s $ p a n a m a b
5 a n a s $ p a n a m a b a n
6 a s $ p a n a m a b a n a n
  b a n a n a s $ p a n a m a
  m a b a n a n a s $ p a n a
  n a m a b a n a n a s $ p a
  n a n a s $ p a n a m a b a
  n a s $ p a n a m a b a n a
  p a n a m a b a n a n a s $
  s $ p a n a m a b a n a n a
```

These strings are sorted

Chop off **a**

```
b a n a n a s $ p a n a m
m a b a n a n a s $ p a n
n a m a b a n a n a s $ p
n a n a s $ p a n a m a b
n a s $ p a n a m a b a n
s $ p a n a m a b a n a n
```

Still
sorted

# Is It True in General?

```
  $panamabananas
1 abananas$panam
2 amabananas$pan
3 anamabananas$p
4 ananas$panamab
5 anas$panamaban
6 as$panamabanan
  bananas$panama
  mabananas$pana
  namabananas$pa
  nanas$panamaba
  nas$panamabana
  panamabananas$
  s$panamabanana
```

These strings are sorted

```
bananas$panam
mabananas$pan
namabananas$p
nanas$panamab
nas$panamaban
s$panamabanan
```

Chop off **a**

Still sorted

Add **a** to end

```
bananas$panama
mabananas$pana
namabananas$pa
nanas$panamaba
nas$panamabana
s$panamabanana
```

# Is It True in General?

```
  $panamabananas
1 abananas$panam
2 amabananas$pan
3 anamabananas$p
4 ananas$panamab
5 anas$panamaban
6 as$panamabanan
  bananas$panama
  mabananas$pana
  namabananas$pa
  nanas$panamaba
  nas$panamabana
  panamabananas$
  s$panamabanana
```

These strings are sorted

Chop off **a**

```
bananas$panam
mabananas$pan
namabananas$p
nanas$panamab
nas$panamaban
s$panamabanan
```

Still sorted

Add **a**
to end

```
bananas$panam**a**
mabananas$pan**a**
namabananas$p**a**
nanas$panamab**a**
nas$panamaban**a**
s$panamabanan**a**
```

Still sorted

# Is It True in General?

```
  $ p a n a m a b a n a n a s
1 a b a n a n a s $ p a n a m
2 a m a b a n a n a s $ p a n
3 a n a m a b a n a n a s $ p
4 a n a n a s $ p a n a m a b
5 a n a s $ p a n a m a b a n
6 a s $ p a n a m a b a n a n
  b a n a n a s $ p a n a m a 1
  m a b a n a n a s $ p a n a 2
  n a m a b a n a n a s $ p a 3
  n a n a s $ p a n a m a b a 4
  n a s $ p a n a m a b a n a 5
  p a n a m a b a n a n a s $
  s $ p a n a m a b a n a n a 6
```

These strings are sorted

Chop off **a**

Ordering doesn't change!

```
b a n a n a s $ p a n a m
m a b a n a n a s $ p a n
n a m a b a n a n a s $ p
n a n a s $ p a n a m a b
n a s $ p a n a m a b a n
s $ p a n a m a b a n a n
```

Still sorted

Add **a** to end

```
b a n a n a s $ p a n a m a
m a b a n a n a s $ p a n a
n a m a b a n a n a s $ p a
n a n a s $ p a n a m a b a
n a s $ p a n a m a b a n a
s $ p a n a m a b a n a n a
```

Still sorted

# First-Last Property

- the *k*-th occurrence of *symbol* in **FirstColumn**

- and the *k*-th occurrence of *symbol* in **LastColumn**

- correspond to appearance of *symbol* at the same position in *Text*.

$p_1a_3n_1a_2m_1a_1b_1a_4n_2a_5n_3a_6s_1\$_1$

$\$_1$panamabanana$s_1$
$a_1$bananas\$panam$_1$
$a_2$mabananas\$pa$n_1$
$a_3$namabananas\$$p_1$
$a_4$nanas\$panama$b_1$
$a_5$nas\$panamaba$n_2$
$a_6$s\$panamabana$n_3$
$b_1$ananas\$panam$a_1$
$m_1$abananas\$pan$a_2$
$n_1$amabananas\$p$a_3$
$n_2$anas\$panamab$a_4$
$n_3$as\$panamaban$a_5$
$p_1$anamabananas$\$_1$
$s_1$\$panamabanan$a_6$

# Inverting BWT Again

$\$_1$panamabananas$_1$
a$_1$bananas\$panam$_1$
a$_2$mabananas\$pan$_1$
a$_3$namabananas\$p$_1$
a$_4$nanas\$panamab$_1$
a$_5$nas\$panamaban$_2$
a$_6$s\$panamabanan$_3$
b$_1$ananas\$panama$_1$
m$_1$abananas\$pana$_2$
n$_1$amabananas\$pa$_3$
n$_2$anas\$panamaba$_4$
n$_3$as\$panamabana$_5$
p$_1$anamabananas$\$_1$
s$_1$\$panamabanana$_6$

# Inverting BWT Again

$\$_1$panamabanana$s_1$
$a_1$bananas\$pana$m_1$
$a_2$mabananas\$pa$n_1$
$a_3$namabananas\$$p_1$
$a_4$nanas\$panama$b_1$
$a_5$nas\$panamaba$n_2$
$a_6$s\$panamabana$n_3$
$b_1$ananas\$panam$a_1$
$m_1$abananas\$pan$a_2$
$n_1$amabananas\$p$a_3$
$n_2$anas\$panamab$a_4$
$n_3$as\$panamaban$a_5$
$p_1$anamabananas$\$_1$
$s_1$\$panamabanan$a_6$

**$\$**   p   a

**s**                    n

              a                                   a

n                                          m

       a                              a

          n              b
                a

# Inverting BWT Again

$$\$_1 \text{panamabanana} \mathbf{s_1}$$
$$a_1 \text{bananas}\$\text{pana} m_1$$
$$a_2 \text{mabananas}\$\text{pa} n_1$$
$$a_3 \text{namabananas}\$ p_1$$
$$a_4 \text{nanas}\$\text{panama} b_1$$
$$a_5 \text{nas}\$\text{panamaban} n_2$$
$$a_6 s\$\text{panamabanan} n_3$$
$$b_1 \text{ananas}\$\text{panam} a_1$$
$$m_1 \text{abananas}\$\text{pan} a_2$$
$$n_1 \text{amabananas}\$\text{p} a_3$$
$$n_2 \text{anas}\$\text{panamab} a_4$$
$$n_3 \text{as}\$\text{panamaban} a_5$$
$$p_1 \text{anamabananas}\$_1$$
$$\mathbf{s_1} \$\text{panamabanan} a_6$$

# Inverting BWT Again

$\$_1$panamabanana$s_1$
$a_1$bananas\$pana$m_1$
$a_2$mabananas\$pa$n_1$
$a_3$namabananas\$$p_1$
$a_4$nanas\$panama$b_1$
$a_5$nas\$panamaba$n_2$
$a_6$s\$panamabana$n_3$
$b_1$ananas\$panam$a_1$
$m_1$abananas\$pan$a_2$
$n_1$amabananas\$p$a_3$
$n_2$anas\$panamab$a_4$
$n_3$as\$panamaban$a_5$
$p_1$anamabananas$\$_1$
**$s_1$**\$panamabanan**$a_6$**

# Inverting BWT Again

$\$_1$panamabanana$s_1$
$a_1$bananas\$pana$m_1$
$a_2$mabananas\$pa$n_1$
$a_3$namabananas\$$p_1$
$a_4$nanas\$panama$b_1$
$a_5$nas\$panamaba$n_2$
$a_6$s\$panamabana$n_3$
$b_1$ananas\$panam$a_1$
$m_1$abananas\$pan$a_2$
$n_1$amabananas\$p$a_3$
$n_2$anas\$panamab$a_4$
$n_3$as\$panamaban$a_5$
$p_1$anamabananas$\$_1$
$s_1$\$panamabanan$a_6$

# Inverting BWT Again

$\$_1$ p a n a m a b a n a n a s $_1$
a $_1$ b a n a n a s \$ p a n a m $_1$
a $_2$ m a b a n a n a s \$ p a n $_1$
a $_3$ n a m a b a n a n a s \$ p $_1$
a $_4$ n a n a s \$ p a n a m a b $_1$
a $_5$ n a s \$ p a n a m a b a n $_2$
**a** $_6$ s \$ p a n a m a b a n a **n** $_3$
b $_1$ a n a n a s \$ p a n a m a $_1$
m $_1$ a b a n a n a s \$ p a n a $_2$
n $_1$ a m a b a n a n a s \$ p a $_3$
n $_2$ a n a s \$ p a n a m a b a $_4$
n $_3$ a s \$ p a n a m a b a n a $_5$
p $_1$ a n a m a b a n a n a s \$ $_1$
s $_1$ \$ p a n a m a b a n a n a $_6$

**a**

**n**

# Inverting BWT Again



$\$_1$ p a n a m a b a n a n a s$_1$
a$_1$ b a n a n a s $ p a n a m$_1$
a$_2$ m a b a n a n a s $ p a n$_1$
a$_3$ n a m a b a n a n a s $ p$_1$
a$_4$ n a n a s $ p a n a m a b$_1$
a$_5$ n a s $ p a n a m a b a n$_2$
a$_6$ s $ p a n a m a b a n a **n$_3$**
b$_1$ a n a n a s $ p a n a m a$_1$
m$_1$ a b a n a n a s $ p a n a$_2$
n$_1$ a m a b a n a n a s $ p a$_3$
n$_2$ a n a s $ p a n a m a b a$_4$
**n$_3$** a s $ p a n a m a b a n a$_5$
p$_1$ a n a m a b a n a n a s $_1$
s$_1$ $ p a n a m a b a n a n a$_6$

# Inverting BWT Again

$$\$_1 p a n a m a b a n a n a s_1$$
$$a_1 b a n a n a s \$ p a n a m_1$$
$$a_2 m a b a n a n a s \$ p a n_1$$
$$a_3 n a m a b a n a n a s \$ p_1$$
$$a_4 n a n a s \$ p a n a m a b_1$$
$$a_5 n a s \$ p a n a m a b a n_2$$
$$a_6 s \$ p a n a m a b a n a n_3$$
$$b_1 a n a n a s \$ p a n a m a_1$$
$$m_1 a b a n a n a s \$ p a n a_2$$
$$n_1 a m a b a n a n a s \$ p a_3$$
$$n_2 a n a s \$ p a n a m a b a_4$$
$$\textbf{n}_3 a s \$ p a n a m a b a n \textbf{a}_5$$
$$p_1 a n a m a b a n a n a s \$_1$$
$$s_1 \$ p a n a m a b a n a n a_6$$

$\$$  p  a

s                n

a                a

**n**            m

**a**            a

n      b

a

# Inverting BWT Again

$\$_1$ p a n a m a b a n a n a s$_1$
a$_1$ b a n a n a s $ p a n a m$_1$
a$_2$ m a b a n a n a s $ p a n$_1$
a$_3$ n a m a b a n a n a s $ p$_1$
a$_4$ n a n a s $ p a n a m a b$_1$
**a$_5$** n a s $ p a n a m a b a n$_2$
a$_6$ s $ p a n a m a b a n a n$_3$
b$_1$ a n a n a s $ p a n a m a$_1$
m$_1$ a b a n a n a s $ p a n a$_2$
n$_1$ a m a b a n a n a s $ p a$_3$
n$_2$ a n a s $ p a n a m a b a$_4$
n$_3$ a s $ p a n a m a b a n **a$_5$**
p$_1$ a n a m a b a n a n a s $_1$
s$_1$ $ p a n a m a b a n a n a$_6$

# Inverting BWT Again

$\$_1$ p a n a m a b a n a n a s $_1$
a $_1$ b a n a n a s \$ p a n a m $_1$
a $_2$ m a b a n a n a s \$ p a n $_1$
a $_3$ n a m a b a n a n a s \$ p $_1$
a $_4$ n a n a s \$ p a n a m a b $_1$
**a $_5$** n a s \$ p a n a m a b a **n $_2$**
a $_6$ s \$ p a n a m a b a n a n $_3$
b $_1$ a n a n a s \$ p a n a m a $_1$
m $_1$ a b a n a n a s \$ p a n a $_2$
n $_1$ a m a b a n a n a s \$ p a $_3$
n $_2$ a n a s \$ p a n a m a b a $_4$
n $_3$ a s \$ p a n a m a b a n a $_5$
p $_1$ a n a m a b a n a n a s \$ $_1$
s $_1$ \$ p a n a m a b a n a n a $_6$

# Inverting BWT Again

$$\$_1 \text{panamabanana} s_1$$
$$a_1 \text{bananas}\$\text{panam}_1$$
$$a_2 \text{mabananas}\$\text{pan}_1$$
$$a_3 \text{namabananas}\$p_1$$
$$a_4 \text{nanas}\$\text{panamab}_1$$
$$a_5 \text{nas}\$\text{panamaba}\textcolor{red}{\mathbf{n_2}}$$
$$a_6 \text{s}\$\text{panamaban}_3$$
$$b_1 \text{ananas}\$\text{panama}_1$$
$$m_1 \text{abananas}\$\text{pan}a_2$$
$$n_1 \text{amabananas}\$\text{p}a_3$$
$$\textcolor{red}{\mathbf{n_2}}\text{anas}\$\text{panamab}a_4$$
$$n_3 \text{as}\$\text{panamaban}a_5$$
$$p_1 \text{anamabananas}\$_1$$
$$s_1 \$\text{panamabanan}a_6$$

# Inverting BWT Again



$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamabanan3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
**n2**anas$panamab**a4**
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6

$_1$ p a n a m a b a n a n a s$_1$
a$_1$ b a n a n a s $ p a n a m$_1$
a$_2$ m a b a n a n a s $ p a n$_1$
a$_3$ n a m a b a n a n a s $ p$_1$
**a**$_4$ n a n a s $ p a n a m a b$_1$
a$_5$ n a s $ p a n a m a b a n$_2$
a$_6$ s $ p a n a m a b a n a n$_3$
b$_1$ a n a n a s $ p a n a m a$_1$
m$_1$ a b a n a n a s $ p a n a$_2$
n$_1$ a m a b a n a n a s $ p a$_3$
n$_2$ a n a s $ p a n a m a b **a**$_4$
n$_3$ a s $ p a n a m a b a n a$_5$
p$_1$ a n a m a b a n a n a s $_1$
s$_1$ $ p a n a m a b a n a n a$_6$

$\$_1$panamabananas$_1$
a$_1$bananas\$panam$_1$
a$_2$mabananas\$pan$_1$
a$_3$namabananas\$p$_1$
**a$_4$**nanas\$panama**b$_1$**
a$_5$nas\$panamaban$_2$
a$_6$s\$panamabanan$_3$
b$_1$ananas\$panama$_1$
m$_1$abananas\$pana$_2$
n$_1$amabananas\$pa$_3$
n$_2$anas\$panamaba$_4$
n$_3$as\$panamabana$_5$
p$_1$anamabananas\$$_1$
s$_1$\$panamabanana$_6$

# Inverting BWT Again

$\$_1$ p a n a m a b a n a n a $s_1$
$a_1$ b a n a n a s $ p a n a $m_1$
$a_2$ m a b a n a n a s $ p a $n_1$
$a_3$ n a m a b a n a n a s $ $p_1$
$a_4$ n a n a s $ p a n a m a $\mathbf{b_1}$
$a_5$ n a s $ p a n a m a b a $n_2$
$a_6$ s $ p a n a m a b a n a $n_3$
$\mathbf{b_1}$ a n a n a s $ p a n a m $a_1$
$m_1$ a b a n a n a s $ p a n $a_2$
$n_1$ a m a b a n a n a s $ p $a_3$
$n_2$ a n a s $ p a n a m a b $a_4$
$n_3$ a s $ p a n a m a b a n $a_5$
$p_1$ a n a m a b a n a n a s $\$_1$
$s_1$ $ p a n a m a b a n a n $a_6$

# Inverting BWT Again

$\$_1$panamabanana$s_1$
$a_1$bananas\$pana$m_1$
$a_2$mabananas\$pa$n_1$
$a_3$namabananas\$$p_1$
$a_4$nanas\$panama$b_1$
$a_5$nas\$panamaba$n_2$
$a_6$s\$panamabana$n_3$
**$b_1$**ananas\$panam**$a_1$**
$m_1$abananas\$pan$a_2$
$n_1$amabananas\$p$a_3$
$n_2$anas\$panamab$a_4$
$n_3$as\$panamaban$a_5$
$p_1$anamabananas$\$_1$
$s_1$\$panamabanan$a_6$

# Inverting BWT Again

$\$_1$ p a n a m a b a n a n a s $s_1$
**$a_1$** b a n a n a s $ p a n a m $m_1$
$a_2$ m a b a n a n a s $ p a n $n_1$
$a_3$ n a m a b a n a n a s $ p $p_1$
$a_4$ n a n a s $ p a n a m a b $b_1$
$a_5$ n a s $ p a n a m a b a n $n_2$
$a_6$ s $ p a n a m a b a n a n $n_3$
$b_1$ a n a n a s $ p a n a m **$a_1$**
$m_1$ a b a n a n a s $ p a n $a_2$
$n_1$ a m a b a n a n a s $ p $a_3$
$n_2$ a n a s $ p a n a m a b $a_4$
$n_3$ a s $ p a n a m a b a n $a_5$
$p_1$ a n a m a b a n a n a s $\$_1$
$s_1$ $ p a n a m a b a n a n $a_6$

# Inverting BWT Again

$\$_1$panamabananas$_1$
**a$_1$**bananas$pana**m$_1$**
a$_2$mabananas$pan$_1$
a$_3$namabananas$p$_1$
a$_4$nanas$panamab$_1$
a$_5$nas$panamaban$_2$
a$_6$s$panamabanan$_3$
b$_1$ananas$panama$_1$
m$_1$abananas$pana$_2$
n$_1$amabananas$pa$_3$
n$_2$anas$panamaba$_4$
n$_3$as$panamabana$_5$
p$_1$anamabananas$\$_1$
s$_1$$panamabanana$_6$

# Inverting BWT Again

# Inverting BWT Again

$\$_1$panamabanana$s_1$
$a_1$bananas\$pana$m_1$
$a_2$mabananas\$pa$n_1$
$a_3$namabananas\$$p_1$
$a_4$nanas\$panama$b_1$
$a_5$nas\$panamaba$n_2$
$a_6$s\$panamabana$n_3$
$b_1$ananas\$panam$a_1$
**$m_1$**abananas\$pan**$a_2$**
$n_1$amabananas\$p$a_3$
$n_2$anas\$panamab$a_4$
$n_3$as\$panamaban$a_5$
$p_1$anamabananas$\$_1$
$s_1$\$panamabanan$a_6$

# Inverting BWT Again

$\$_1$panamabananas$_1$
a$_1$bananas$panam$_1$
**a$_2$**mabananas$pan$_1$
a$_3$namabananas$p$_1$
a$_4$nanas$panamab$_1$
a$_5$nas$panamaban$_2$
a$_6$s$panamabanan$_3$
b$_1$ananas$panama$_1$
m$_1$abananas$pan**a$_2$**
n$_1$amabananas$pa$_3$
n$_2$anas$panamaba$_4$
n$_3$as$panamaban a$_5$
p$_1$anamabananas$\$_1$
s$_1$$panamabanan a$_6$

# Inverting BWT Again

$\$_1$panamabananas$_1$
a$_1$bananas\$panam$_1$
**a$_2$**mabananas\$pa**n$_1$**
a$_3$namabananas\$p$_1$
a$_4$nanas\$panamab$_1$
a$_5$nas\$panaman$_2$
a$_6$s\$panamabanan$_3$
b$_1$ananas\$panama$_1$
m$_1$abananas\$pana$_2$
n$_1$amabananas\$pa$_3$
n$_2$anas\$panamaba$_4$
n$_3$as\$panamabana$_5$
p$_1$anamabananas$\$_1$
s$_1$\$panamabanana$_6$

# Inverting BWT Again

$\$_1$panamabanana$s_1$
$a_1$bananas\$pana$m_1$
$a_2$mabananas\$pa$\mathbf{n_1}$
$a_3$namabananas\$$p_1$
$a_4$nanas\$panama$b_1$
$a_5$nas\$panamaba$n_2$
$a_6$s\$panamabana$n_3$
$b_1$ananas\$panam$a_1$
$m_1$abananas\$pan$a_2$
$\mathbf{n_1}$amabananas\$p$a_3$
$n_2$anas\$panamab$a_4$
$n_3$as\$panamaban$a_5$
$p_1$anamabananas$\$_1$
$s_1$\$panamabanan$a_6$

# Inverting BWT Again



$_1$ p a n a m a b a n a n a s $_1$
a $_1$ b a n a n a s $ p a n a m $_1$
a $_2$ m a b a n a n a s $ p a n $_1$
a $_3$ n a m a b a n a n a s $ p $_1$
a $_4$ n a n a s $ p a n a m a b $_1$
a $_5$ n a s $ p a n a m a b a n $_2$
a $_6$ s $ p a n a m a b a n a n $_3$
b $_1$ a n a n a s $ p a n a m a $_1$
m $_1$ a b a n a n a s $ p a n a $_2$
**n $_1$** a m a b a n a n a s $ p **a $_3$**
n $_2$ a n a s $ p a n a m a b a $_4$
n $_3$ a s $ p a n a m a b a n a $_5$
p $_1$ a n a m a b a n a n a s $_1$
s $_1$ $ p a n a m a b a n a n a $_6$

# Inverting BWT Again

$\$_1$ p a n a m a b a n a n a $s_1$
$a_1$ b a n a n a s \$ p a n a $m_1$
$a_2$ m a b a n a n a s \$ p a $n_1$
$\mathbf{a_3}$ n a m a b a n a n a s \$ $p_1$
$a_4$ n a n a s \$ p a n a m a $b_1$
$a_5$ n a s \$ p a n a m a b a $n_2$
$a_6$ s \$ p a n a m a b a n a $n_3$
$b_1$ a n a n a s \$ p a n a m $a_1$
$m_1$ a b a n a n a s \$ p a n $a_2$
$n_1$ a m a b a n a n a s \$ p $\mathbf{a_3}$
$n_2$ a n a s \$ p a n a m a b $a_4$
$n_3$ a s \$ p a n a m a b a n $a_5$
$p_1$ a n a m a b a n a n a s $\$_1$
$s_1$ \$ p a n a m a b a n a n $a_6$

# Inverting BWT Again

$\$_1$ p a n a m a b a n a n a s $_1$
a $_1$ b a n a n a s $ p a n a m $_1$
a $_2$ m a b a n a n a s $ p a n $_1$
**a $_3$** n a m a b a n a n a s $ **p $_1$**
a $_4$ n a n a s $ p a n a m a b $_1$
a $_5$ n a s $ p a n a m a b a n $_2$
a $_6$ s $ p a n a m a b a n a n $_3$
b $_1$ a n a n a s $ p a n a m a $_1$
m $_1$ a b a n a n a s $ p a n a $_2$
n $_1$ a m a b a n a n a s $ p a $_3$
n $_2$ a n a s $ p a n a m a b a $_4$
n $_3$ a s $ p a n a m a b a n a $_5$
p $_1$ a n a m a b a n a n a s $_1$
s $_1$ $ p a n a m a b a n a n a $_6$

# Inverting BWT Again

$\$_1$panamabanana$s_1$
$a_1$bananas\$pana$m_1$
$a_2$mabananas\$pa$n_1$
$a_3$namabananas\$**p₁**
$a_4$nanas\$panama$b_1$
$a_5$nas\$panamaba$n_2$
$a_6$s\$panamabana$n_3$
$b_1$ananas\$panam$a_1$
$m_1$abananas\$pan$a_2$
$n_1$amabananas\$p$a_3$
$n_2$anas\$panamab$a_4$
$n_3$as\$panamaban$a_5$
**p₁**anamabananas$\$_1$
$s_1$\$panamabanan$a_6$

# We Are Done!

$\$_1$panamabananas$_1$
a$_1$bananas\$pana m$_1$
a$_2$mabananas\$pa n$_1$
a$_3$namabananas\$ p$_1$
a$_4$nanas\$panama b$_1$
a$_5$nas\$panamaba n$_2$
a$_6$s\$panamabana n$_3$
b$_1$ananas\$panam a$_1$
m$_1$abananas\$pan a$_2$
n$_1$amabananas\$p a$_3$
n$_2$anas\$panamab a$_4$
n$_3$as\$panamaban a$_5$
p$_1$anamabananas $\$_1$
s$_1$\$panamabanan a$_6$

# This Was Fast!



- Memory: 2|*Text*|
- Time: O(|*Text*|)

# Outline

- Burrows-Wheeler Transform

- Inverting Burrows-Wheeler Transform

- **Using BWT for Pattern Matching**

- Suffix Arrays

- Approximate Pattern Matching

# Back to Pattern Matching

- Suffix Tree Pattern Matching:
  - Runtime: $O(|Text| + |Patterns|)$
  - Memory: $20 \cdot |Text|$

For human genome:



- $|Text| \approx 3*10^9$

- Can we use BWT($Text$) to design a more memory efficient linear-time algorithm for Multiple Pattern Matching?

# Finding Pattern Matches Using BWT

- Searching for **ana** in p**ana**mab**anana**s

$$\$_1 \text{panamabanana} s_1$$
$$a_1 \text{bananas\$pana} m_1$$
$$a_2 \text{mabananas\$pa} n_1$$
$$\mathbf{a_3 na} \text{mabananas\$p} 1$$
$$\mathbf{a_4 na} \text{nas\$panama} b_1$$
$$\mathbf{a_5 na} \text{s\$panamaba} n_2$$
$$a_6 \text{s\$panamabana} n_3$$
$$b_1 \text{ananas\$panama} a_1$$
$$m_1 \text{abananas\$pan} a_2$$
$$n_1 \text{amabananas\$p} a_3$$
$$n_2 \text{anas\$panamab} a_4$$
$$n_3 \text{as\$panamaban} a_5$$
$$p_1 \text{anamabananas} \$_1$$
$$s_1 \text{\$panamabanan} a_6$$

# Lets Start by Matching the Last Symbol (a)

- Searching for an**a** in panamabananas

$$\$_1 \text{panamabanana} s_1$$
$$\mathbf{a}_1 \text{bananas\$pana} m_1$$
$$\mathbf{a}_2 \text{mabananas\$pa} n_1$$
$$\mathbf{a}_3 \text{namabananas\$} p_1$$
$$\mathbf{a}_4 \text{nanas\$panama} b_1$$
$$\mathbf{a}_5 \text{nas\$panamaban} n_2$$
$$\mathbf{a}_6 \text{s\$panamabanan} n_3$$
$$b_1 \text{ananas\$panam} a_1$$
$$m_1 \text{abananas\$pan} a_2$$
$$n_1 \text{amabananas\$p} a_3$$
$$n_2 \text{anas\$panamab} a_4$$
$$n_3 \text{as\$panamaban} a_5$$
$$p_1 \text{anamabananas\$}_1$$
$$s_1 \text{\$panamabanan} a_6$$

# Matching the Last Two Symbols (na)

- Searching for a**na** in panamabananas

$\$_1$ p a n a m a b a n a n a $s_1$
**a**$_1$ b a n a n a s \$ p a n a **m**$_1$
**a**$_2$ m a b a n a n a s \$ p a **n**$_1$
**a**$_3$ n a m a b a n a n a s \$ **p**$_1$
**a**$_4$ n a n a s \$ p a n a m a **b**$_1$
**a**$_5$ n a s \$ p a n a m a b a **n**$_2$
**a**$_6$ s \$ p a n a m a b a n a **n**$_3$
b$_1$ a n a n a s \$ p a n a m a$_1$
m$_1$ a b a n a n a s \$ p a n a$_2$
n$_1$ a m a b a n a n a s \$ p a$_3$
n$_2$ a n a s \$ p a n a m a b a$_4$
n$_3$ a s \$ p a n a m a b a n a$_5$
p$_1$ a n a m a b a n a n a s \$$_1$
s$_1$ \$ p a n a m a b a n a n a$_6$

# Three Matches of na Found!

- Searching for a**na** in panamabananas

```
$₁ p a n a m a b a n a n a s₁
a₁ b a n a n a s $ p a n a m₁
a₂ m a b a n a n a s $ p a n₁
a₃ n a m a b a n a n a s $ p₁
a₄ n a n a s $ p a n a m a b₁
a₅ n a s $ p a n a m a b a n₂
a₆ s $ p a n a m a b a n a n₃
b₁ a n a n a s $ p a n a m a₁
m₁ a b a n a n a s $ p a n a₂
n₁ a m a b a n a n a s $ p a₃
n₂ a n a s $ p a n a m a b a₄
n₃ a s $ p a n a m a b a n a₅
p₁ a n a m a b a n a n a s $₁
s₁ $ p a n a m a b a n a n a₆
```

# Three Matches of na Found!

- Searching for a**na** in panamabananas

$\$_1$ panamabanana $s_1$
$a_1$ bananas\$pana $m_1$
$a_2$ mabananas\$pa **$n_1$**
$a_3$ namabananas\$p $p_1$
$a_4$ nanas\$panama $b_1$
$a_5$ nas\$panamaba **$n_2$**
$a_6$ s\$panamabana **$n_3$**
$b_1$ ananas\$panam $a_1$
$m_1$ abananas\$pan $a_2$
**$n_1$** amabananas\$p $a_3$
**$n_2$** anas\$panamab $a_4$
**$n_3$** as\$panamaban $a_5$
$p_1$ anamabananas\$ $\$_1$
$s_1$ \$panamabanan $a_6$

# Three Matches of na Found!

- Searching for a**na** in panamabananas

$\$_1$ p a n a m a b a n a n a s$_1$
a$_1$ b a n a n a s \$ p a n a m$_1$
a$_2$ m a b a n a n a s \$ p a n$_1$
a$_3$ n a m a b a n a n a s \$ p$_1$
a$_4$ n a n a s \$ p a n a m a b$_1$
a$_5$ n a s \$ p a n a m a b a n$_2$
a$_6$ s \$ p a n a m a b a n a n$_3$
b$_1$ a n a n a s \$ p a n a m a$_1$
m$_1$ a b a n a n a s \$ p a n a$_2$
**n$_1$ a** m a b a n a n a s \$ p a$_3$
**n$_2$ a** n a s \$ p a n a m a b a$_4$
**n$_3$ a** s \$ p a n a m a b a n a$_5$
p$_1$ a n a m a b a n a n a s \$$_1$
s$_1$ \$ p a n a m a b a n a n a$_6$

# Matching **ana**

- Searching for **ana** in `panamabananas`

```
$₁ p a n a m a b a n a n a s₁
a₁ b a n a n a s $ p a n a m₁
a₂ m a b a n a n a s $ p a n₁
a₃ n a m a b a n a n a s $ p₁
a₄ n a n a s $ p a n a m a b₁
a₅ n a s $ p a n a m a b a n₂
a₆ s $ p a n a m a b a n a n₃
b₁ a n a n a s $ p a n a m a₁
m₁ a b a n a n a s $ p a n a₂
n₁ a m a b a n a n a s $ p a₃
n₂ a n a s $ p a n a m a b a₄
n₃ a s $ p a n a m a b a n a₅
p₁ a n a m a b a n a n a s $₁
s₁ $ p a n a m a b a n a n a₆
```

# Three Matches of ana Found!

- Searching for **ana** in panamabananas

$\$_1$panamabana s$_1$
a$_1$bananas$\$$pana m$_1$
a$_2$mabananas$\$$pa n$_1$
**a$_3$na**mabananas$\$$p$_1$
**a$_4$na**nas$\$$panama b$_1$
**a$_5$na**s$\$$panamaba n$_2$
a$_6$s$\$$panamabana n$_3$
b$_1$ananas$\$$panam a$_1$
m$_1$abananas$\$$pan a$_2$
n$_1$amabananas$\$$p a$_3$
n$_2$anas$\$$panamab a$_4$
n$_3$as$\$$panamaban a$_5$
p$_1$anamabananas$\$_1$
s$_1$$\$$panamabanan a$_6$

# Searching for **ana** using *top* and *bottom* pointers



topIndex ← first position of *symbol* among positions from *top* to *bottom* in *LastColumn*

bottomIndex ← last position of *symbol* among positions from *top* to *bottom* in *LastColumn*

# BWMatching

**BWMATCHING**(*FirstColumn, LastColumn, Pattern,* LASTTOFIRST)
  *top* ← 0
  *bottom* ← |*LastColumn*| − 1
  **while** *top* ≤ *bottom*
      **if** *Pattern* is nonempty
          *symbol* ← last letter in *Pattern*
          remove last letter from *Pattern*
          **if** positions from *top* to *bottom* in *LastColumn* contain *symbol*
              *topIndex* ← first position of *symbol* among positions from *top* to *bottom*
                               in *LastColumn*
              *bottomIndex* ← last position of *symbol* among positions from *top* to
                               *bottom* in *LastColumn*
              *top* ← LASTTOFIRST(*topIndex*)
              *bottom* ← LASTTOFIRST(*bottomIndex*)
          **else**
              **return** 0
      **else**
          **return** *bottom* − *top* + 1

Given a symbol at position *index* in *LastColumn,*
**LastToFirst**(*index*) defines the position of this symbol in *FirstColumn*

# BWMatching is slow:
## it analyzes every symbol from *top* to *bottom* in each step!



**if** positions from *top* to *bottom* in *LastColumn* contain *symbol*

    *topIndex* ← first position of *symbol* among positions from *top* to *bottom* in *LastColumn*

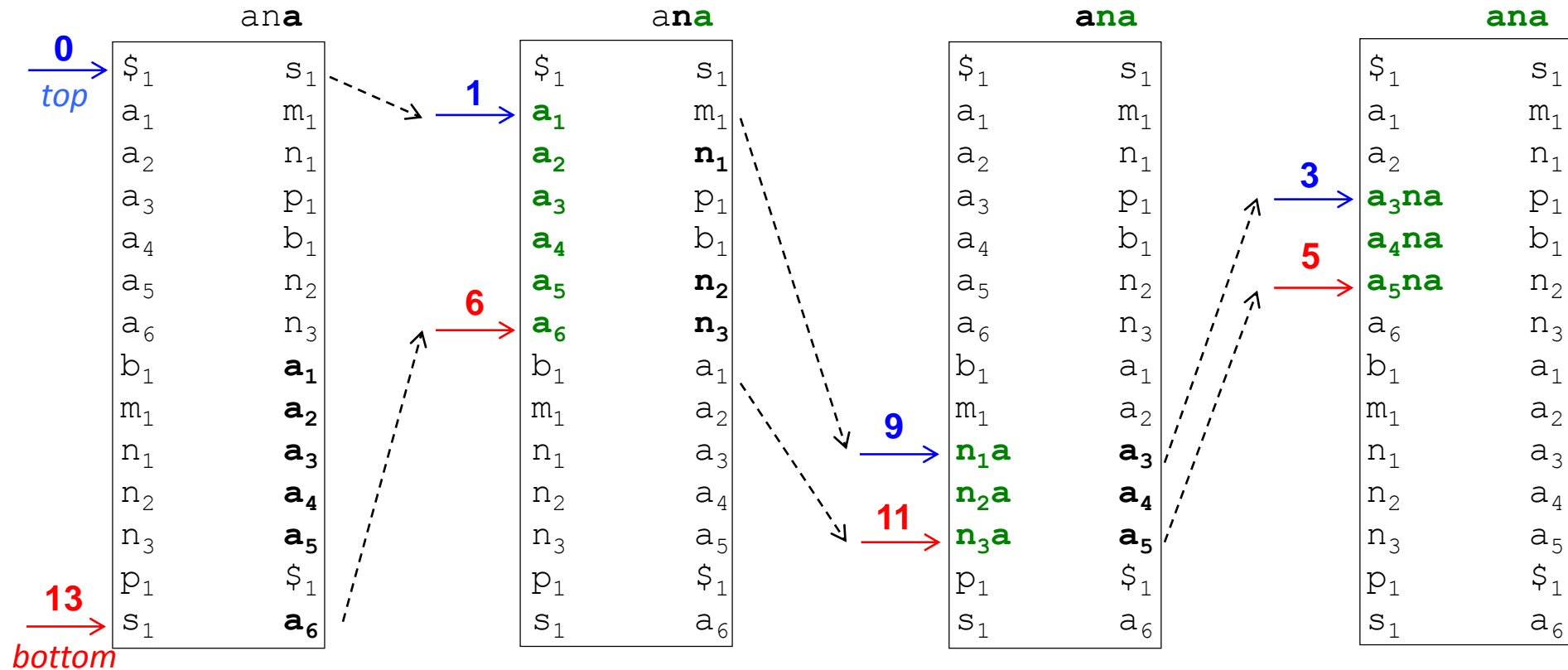    *bottomIndex* ← last position of *symbol* among positions from *top* to *bottom* in *LastColumn*

# Introducing *Count* Array

| i | FirstColumn | LastColumn | LastToFirst(i) | $ | a | b | m | n | p | s |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | $_1$ | s$_1$ | 13 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| 1 | a$_1$ | m$_1$ | 8 | 0 | 0 | 0 | 1 | 0 | 0 | 1 |
| 2 | a$_2$ | n$_1$ | 9 | 0 | 0 | 0 | 1 | 1 | 0 | 1 |
| 3 | a$_3$ | p$_1$ | 12 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| 4 | a$_4$ | b$_1$ | 7 | 0 | 0 | 1 | 1 | 1 | 1 | 1 |
| 5 | a$_5$ | n$_2$ | 10 | 0 | 0 | 1 | 1 | 2 | 1 | 1 |
| 6 | a$_6$ | n$_3$ | 11 | 0 | 0 | 1 | 1 | 3 | 1 | 1 |
| 7 | b$_1$ | a$_1$ | 1 | 0 | 1 | 1 | 1 | 3 | 1 | 1 |
| 8 | m$_1$ | a$_2$ | 2 | 0 | 2 | 1 | 1 | 3 | 1 | 1 |
| 9 | n$_1$ | a$_3$ | 3 | 0 | 3 | 1 | 1 | 3 | 1 | 1 |
| 10 | n$_2$ | a$_4$ | 4 | 0 | 4 | 1 | 1 | 3 | 1 | 1 |
| 11 | n$_3$ | a$_5$ | 5 | 0 | 5 | 1 | 1 | 3 | 1 | 1 |
| 12 | p$_1$ | $_1$ | 0 | 1 | 5 | 1 | 1 | 3 | 1 | 1 |
| 13 | s$_1$ | a$_6$ | 6 | 1 | 6 | 1 | 1 | 3 | 1 | 1 |

Count$_{symbol}$(i, LastColumn):
#occurrences of *symbol* in the first *i* positions of *LastColumn*

# BetterBWMatching

```
BETTERBWMATCHING(FIRSTOCCURRENCE, LastColumn, Pattern, COUNT)
    top ← 0
    bottom ← |LastColumn| − 1
    while top ≤ bottom
        if Pattern is nonempty
            symbol ← last letter in Pattern
            remove last letter from Pattern
            top ← FIRSTOCCURRENCE(symbol) + COUNT_symbol(top, LastColumn)
            bottom ← FIRSTOCCURRENCE(symbol) + COUNT_symbol(bottom + 1,
                LastColumn) − 1
        else
            return bottom − top + 1
    return
```

BIOINFORMATICS ALGORITHMS
An Active Learning Approach

by Phillip Compeau & Pavel Pevzner

# Where Are the Matches?

- We know that **ana** occurs 3 times, but where does **ana** appear in *Text*???

$\$_1$panamabanana$s_1$
$a_1$bananas\$panam$_1$
$a_2$mabananas\$pan$_1$
**$a_3$na**mabananas\$p$_1$
**$a_4$na**nas\$panamab$_1$
**$a_5$na**s\$panamaban$_2$
$a_6$s\$panamabanan$_3$
$b_1$ananas\$panama$_1$
$m_1$abananas\$pana$_2$
$n_1$amabananas\$pa$_3$
$n_2$anas\$panamaba$_4$
$n_3$as\$panamabana$_5$
$p_1$anamabananas$\$_1$
$s_1\$$panamabanana$_6$

# Outline

- [Burrows-Wheeler Transform](#)

- [Inverting Burrows-Wheeler Transform](#)

- [Using BWT for Pattern Matching](#)

- **[Suffix Arrays](#)**

- [Approximate Pattern Matching](#)

# Where Are the Matches?

- **Suffix array** holds starting position of each suffix

$\$_1$panamabanana$s_1$
$a_1$bananas\$pana$m_1$
$a_2$mabananas\$pa$n_1$
$a_3$namabananas\$$p_1$
$a_4$nanas\$panama$b_1$
$a_5$nas\$panamaba$n_2$
$a_6$s\$panamabana$n_3$
$b_1$ananas\$panam$a_1$
$m_1$abananas\$pan$a_2$
$n_1$amabananas\$p$a_3$
$n_2$anas\$panamab$a_4$
$n_3$as\$panamaban$a_5$
$p_1$anamabananas$\$_1$
$s_1$\$panamabanan$a_6$

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

$$\texttt{panamabananas\$}$$

| 1 3 |
|---|

$\$_1$panamabananas$_1$
$a_1$bananas\$panam$_1$
$a_2$mabananas\$pan$_1$
$a_3$namabananas\$p$_1$
$a_4$nanas\$panamab$_1$
$a_5$nas\$panamaban$_2$
$a_6$s\$panamabanan$_3$
$b_1$ananas\$panama$_1$
$m_1$abananas\$pana$_2$
$n_1$amabananas\$pa$_3$
$n_2$anas\$panamaba$_4$
$n_3$as\$panamabana$_5$
$p_1$anamabananas\$$_1$
$s_1$\$panamabanana$_6$

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

$panam$**abananas$**

| 13 |
|----|
| 5  |

$\$_1$ p a n a m a b a n a n a s $_1$
$a_1$ **b a n a n a s $** p a n a m $_1$
$a_2$ m a b a n a n a s $ p a n $_1$
$a_3$ n a m a b a n a n a s $ p $_1$
$a_4$ n a n a s $ p a n a m a b $_1$
$a_5$ n a s $ p a n a m a b a n $_2$
$a_6$ s $ p a n a m a b a n $_3$
$b_1$ a n a n a s $ p a n a m a $_1$
$m_1$ a b a n a n a s $ p a n a $_2$
$n_1$ a m a b a n a n a s $ p a $_3$
$n_2$ a n a s $ p a n a m a b a $_4$
$n_3$ a s $ p a n a m a b a n a $_5$
$p_1$ a n a m a b a n a n a s $ $_1$
$s_1$ $ p a n a m a b a n a n a $_6$

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

pan**amabananas$**

```
13
 5
 3
```

$\$_1$ panamabananas$_1$
$a_1$**bananas\$** panam$_1$
$a_2$**mabananas\$** pan$_1$
$a_3$namabananas\$p$_1$
$a_4$nanas\$panamab$_1$
$a_5$nas\$panamaban$_2$
$a_6$s\$panamabanan$_3$
$b_1$ananas\$panama$_1$
$m_1$abananas\$pana$_2$
$n_1$amabananas\$pa$_3$
$n_2$anas\$panamaba$_4$
$n_3$as\$panamabana$_5$
$p_1$anamabananas\$$_1$
$s_1$\$panamabanana$_6$

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

$\text{p}$**anamabananas\$**

| | |
|---|---|
| 13 | $\$_1$panamabananas$_1$ |
| 5 | $a_1$**bananas\$**panam$_1$ |
| 3 | $a_2$**mabananas\$**pan$_1$ |
| 1 | $a_3$**namabananas\$**p$_1$ |
| | $a_4$nanas\$panamab$_1$ |
| | $a_5$nas\$panamaban$_2$ |
| | $a_6$s\$panamabanan$_3$ |
| | $b_1$ananas\$panama$_1$ |
| | $m_1$abananas\$pana$_2$ |
| | $n_1$amabananas\$pa$_3$ |
| | $n_2$anas\$panamaba$_4$ |
| | $n_3$as\$panamabana$_5$ |
| | $p_1$anamabananas\$$_1$ |
| | $s_1$\$panamabanana$_6$ |

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

panamab**ananas$**

| 13 |
|----|
| 5  |
| 3  |
| 1  |
| 7  |

$\$_1$panamabananas$_1$
$a_1$**bananas\$**panam$_1$
$a_2$**mabananas\$**pan$_1$
$a_3$**namabananas\$**p$_1$
$a_4$**nanas\$**panamab$_1$
$a_5$nas\$panamaban$_2$
$a_6$s\$panamabanan$_3$
$b_1$ananas\$panama$_1$
$m_1$abananas\$pana$_2$
$n_1$amabananas\$pa$_3$
$n_2$anas\$panamaba$_4$
$n_3$as\$panamabana$_5$
$p_1$anamabananas\$$_1$
$s_1$\$panamabanana$_6$

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

panamaban**anas$**

| | |
|---|---|
| 13 | $\$_1$panamabananas$_1$ |
| 5 | $a_1$**bananas$**panam$_1$ |
| 3 | $a_2$**mabananas$**pan$_1$ |
| 1 | $a_3$**namabananas$**p$_1$ |
| 7 | $a_4$**nanas$**panamab$_1$ |
| 9 | $a_5$**nas$**panamaban$_2$ |
| | $a_6$s$panamaban$_3$ |
| | $b_1$ananas$panama$_1$ |
| | $m_1$abananas$pana$_2$ |
| | $n_1$amabananas$pa$_3$ |
| | $n_2$anas$panamaba$_4$ |
| | $n_3$as$panamabana$_5$ |
| | $p_1$anamabananas$$_1$ |
| | $s_1$$panamabanana$_6$ |

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

p a n a m a b a n a n**a s $**

| | |
|---|---|
| 13 | $\$_1$ p a n a m a b a n a n a s $_1$ |
| 5 | $a_1$ **b a n a n a s $** p a n a m $_1$ |
| 3 | $a_2$ **m a b a n a n a s $** p a n $_1$ |
| 1 | $a_3$ **n a m a b a n a n a s $** p $_1$ |
| 7 | $a_4$ **n a n a s $** p a n a m a b $_1$ |
| 9 | $a_5$ **n a s $** p a n a m a b a n $_2$ |
| 11 | $a_6$ **s $** p a n a m a b a n a n $_3$ |
| | $b_1$ a n a n a s $ p a n a m a $_1$ |
| | $m_1$ a b a n a n a s $ p a n a $_2$ |
| | $n_1$ a m a b a n a n a s $ p a $_3$ |
| | $n_2$ a n a s $ p a n a m a b a $_4$ |
| | $n_3$ a s $ p a n a m a b a n a $_5$ |
| | $p_1$ a n a m a b a n a n a s $ $_1$ |
| | $s_1$ $ p a n a m a b a n a n a $_6$ |

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

panama**bananas$**

| | |
|---|---|
| 13 | $\$_1$panamabananas$_1$ |
| 5 | $a_1$**bananas\$**panam$_1$ |
| 3 | $a_2$**mabananas\$**pan$_1$ |
| 1 | $a_3$**namabananas\$**p$_1$ |
| 7 | $a_4$**nanas\$**panamab$_1$ |
| 9 | $a_5$**nas\$**panamaban$_2$ |
| 11 | $a_6$**s\$**panamaban$_3$ |
| 6 | $b_1$**ananas\$**panama$_1$ |
| | $m_1$abananas\$pana$_2$ |
| | $n_1$amabananas\$pa$_3$ |
| | $n_2$anas\$panamaba$_4$ |
| | $n_3$as\$panamabana$_5$ |
| | $p_1$anamabananas\$$_1$ |
| | $s_1$\$panamabanana$_6$ |

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

pana**mabananas\$**

| | |
|---|---|
| 13 | **\$$_1$panamabananas$_1$** |
| 5 | **a$_1$bananas\$**panam$_1$ |
| 3 | **a$_2$mabananas\$**pan$_1$ |
| 1 | **a$_3$namabananas\$**p$_1$ |
| 7 | **a$_4$nanas\$**panamab$_1$ |
| 9 | **a$_5$nas\$**panamaban$_2$ |
| 11 | **a$_6$s\$**panamaban$_3$ |
| 6 | **b$_1$ananas\$**panama$_1$ |
| 4 | **m$_1$abananas\$**pana$_2$ |
| | n$_1$amabananas\$pa$_3$ |
| | n$_2$anas\$panamaba$_4$ |
| | n$_3$as\$panamabana$_5$ |
| | p$_1$anamabananas\$$_1$ |
| | s$_1$\$panamabanana$_6$ |

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

$\mathtt{pa}$**namabananas\$**

| | |
|---|---|
| 13 | $\mathbf{\$_1}$panamabananas$_1$ |
| 5 | $\mathbf{a_1 bananas\$}$panam$_1$ |
| 3 | $\mathbf{a_2 mabananas\$}$pan$_1$ |
| 1 | $\mathbf{a_3 namabananas\$}$p$_1$ |
| 7 | $\mathbf{a_4 nanas\$}$panamab$_1$ |
| 9 | $\mathbf{a_5 nas\$}$panamaban$_2$ |
| 11 | $\mathbf{a_6 s\$}$panamabanan$_3$ |
| 6 | $\mathbf{b_1 ananas\$}$panama$_1$ |
| 4 | $\mathbf{m_1 abananas\$}$pana$_2$ |
| 2 | $\mathbf{n_1 amabananas\$}$pa$_3$ |
| | n$_2$anas\$panamaba$_4$ |
| | n$_3$as\$panamabana$_5$ |
| | p$_1$anamabananas\$$_1$ |
| | s$_1$\$panamabanana$_6$ |

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

panamaba**nanas$**

| | |
|---|---|
| 13 | **$**$_1$panamabananas$_1$ |
| 5 | **a**$_1$**bananas$**panam$_1$ |
| 3 | **a**$_2$**mabananas$**pan$_1$ |
| 1 | **a**$_3$**namabananas$**p$_1$ |
| 7 | **a**$_4$**nanas$**panamab$_1$ |
| 9 | **a**$_5$**nas$**panamaban$_2$ |
| 11 | **a**$_6$**s$**panamabanan$_3$ |
| 6 | **b**$_1$**ananas$**panama$_1$ |
| 4 | **m**$_1$**abananas$**pana$_2$ |
| 2 | **n**$_1$**amabananas$**pa$_3$ |
| 8 | **n**$_2$**anas$**panamaba$_4$ |
| | n$_3$as$panamabana$_5$ |
| | p$_1$anamabananas$_1$ |
| | s$_1$$panamabanana$_6$ |

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

panamabana**nas$**

| | |
|---|---|
| 13 | $\$_1$panamabananas$_1$ |
| 5 | $a_1$**bananas\$**panam$_1$ |
| 3 | $a_2$**mabananas\$**pan$_1$ |
| 1 | $a_3$**namabananas\$**p$_1$ |
| 7 | $a_4$**nanas\$**panamab$_1$ |
| 9 | $a_5$**nas\$**panamaban$_2$ |
| 11 | $a_6$**s\$**panamaban$_3$ |
| 6 | $b_1$**ananas\$**panama$_1$ |
| 4 | $m_1$**abananas\$**pana$_2$ |
| 2 | $n_1$**amabananas\$**pa$_3$ |
| 8 | $n_2$**anas\$**panamaba$_4$ |
| 10 | $n_3$**as\$**panamabana$_5$ |
| | $p_1$anamabananas$\$_1$ |
| | $s_1$\$panamabanana$_6$ |

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

**panamabananas$**

| | |
|---|---|
| 13 | $\$_1$panamabananas$_1$ |
| 5 | $a_1$bananas$\$$panam$_1$ |
| 3 | $a_2$mabananas$\$$pan$_1$ |
| 1 | $a_3$namabananas$\$$p$_1$ |
| 7 | $a_4$nanas$\$$panamab$_1$ |
| 9 | $a_5$nas$\$$panamaban$_2$ |
| 11 | $a_6$s$\$$panamabanan$_3$ |
| 6 | $b_1$ananas$\$$panama$_1$ |
| 4 | $m_1$abananas$\$$pana$_2$ |
| 2 | $n_1$amabananas$\$$pa$_3$ |
| 8 | $n_2$anas$\$$panamaba$_4$ |
| 10 | $n_3$as$\$$panamabana$_5$ |
| 0 | $p_1$anamabananas$\$_1$ |
| | $s_1$$\$$panamabanana$_6$ |

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

panamabanana**s$**

| | |
|---|---|
| 13 | $\$_1$panamabananas$_1$ |
| 5 | $a_1$bananas$\$$panam$_1$ |
| 3 | $a_2$mabananas$\$$pan$_1$ |
| 1 | $a_3$namabananas$\$$p$_1$ |
| 7 | $a_4$nanas$\$$panamab$_1$ |
| 9 | $a_5$nas$\$$panamaban$_2$ |
| 11 | $a_6$s$\$$panamaban$_3$ |
| 6 | $b_1$ananas$\$$panama$_1$ |
| 4 | $m_1$abananas$\$$pana$_2$ |
| 2 | $n_1$amabananas$\$$pa$_3$ |
| 8 | $n_2$anas$\$$panamaba$_4$ |
| 10 | $n_3$as$\$$panamabana$_5$ |
| 0 | $p_1$anamabananas$\$_1$ |
| 12 | $s_1$$\$$panamabanana$_6$ |

# Suffix Array

- **Suffix array:** holds starting position of each suffix beginning a row.

| | |
|---|---|
| 13 | $\$_1$panamabananas$_1$ |
| 5 | a$_1$bananas$panam$_1$ |
| 3 | a$_2$mabananas$pan$_1$ |
| 1 | a$_3$namabananas$p$_1$ |
| 7 | a$_4$nanas$panamab$_1$ |
| 9 | a$_5$nas$panamaban$_2$ |
| 11 | a$_6$s$panamabanan$_3$ |
| 6 | b$_1$ananas$panama$_1$ |
| 4 | m$_1$abananas$pana$_2$ |
| 2 | n$_1$amabananas$pa$_3$ |
| 8 | n$_2$anas$panamaba$_4$ |
| 10 | n$_3$as$panamaban$a$_5$ |
| 0 | p$_1$anamabananas$\$_1$ |
| 12 | s$_1$\$panamabanana$_6$ |

# Using the Suffix Array to Find Matches

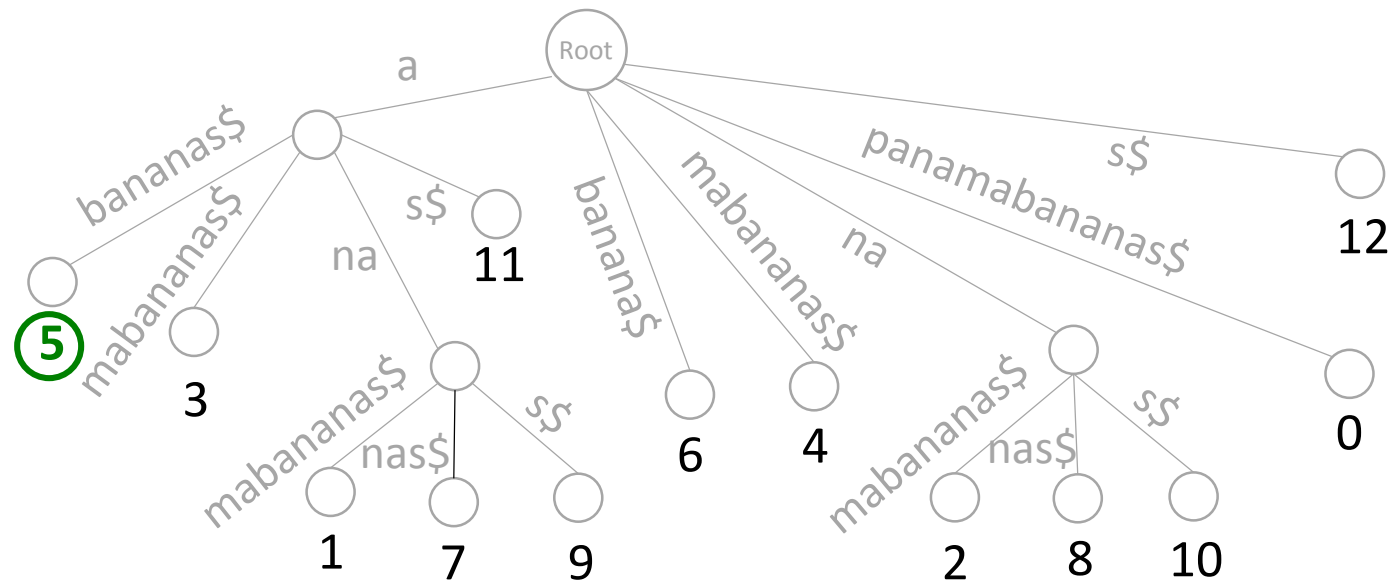- Thus, **ana** occurs at positions **1**, **7**, **9**:

- 

  p**ana**mab**anana**s$

  ↑     ↑ ↑

| | |
|---|---|
| 13 | $_1$panamabanana$s_1$ |
| 5 | $a_1$bananas\$pana$m_1$ |
| 3 | $a_2$mabananas\$pa$n_1$ |
| **1** | **$a_3$na**mabananas\$$p_1$ |
| **7** | **$a_4$na**nas\$panama$b_1$ |
| **9** | **$a_5$na**s\$panamaba$n_2$ |
| 11 | $a_6$s\$panamabana$n_3$ |
| 6 | $b_1$ananas\$panam$a_1$ |
| 4 | $m_1$abananas\$pan$a_2$ |
| 2 | $n_1$amabananas\$p$a_3$ |
| 8 | $n_2$anas\$panamab$a_4$ |
| 10 | $n_3$as\$panamaban$a_5$ |
| 0 | $p_1$anamabananas$\$_1$ |
| 12 | $s_1$\$panamabanan$a_6$ |

**Naïve algorithm for constructing suffix array** (sorting all suffixes of *Text*)
$O(|Text| \cdot \log|Text|)$ comparisons

# From Suffix Tree to Suffix Array:
# Depth-First Traversal



[13  (5)  3  1  7  9  11  6  4  2  8  10  0  12]
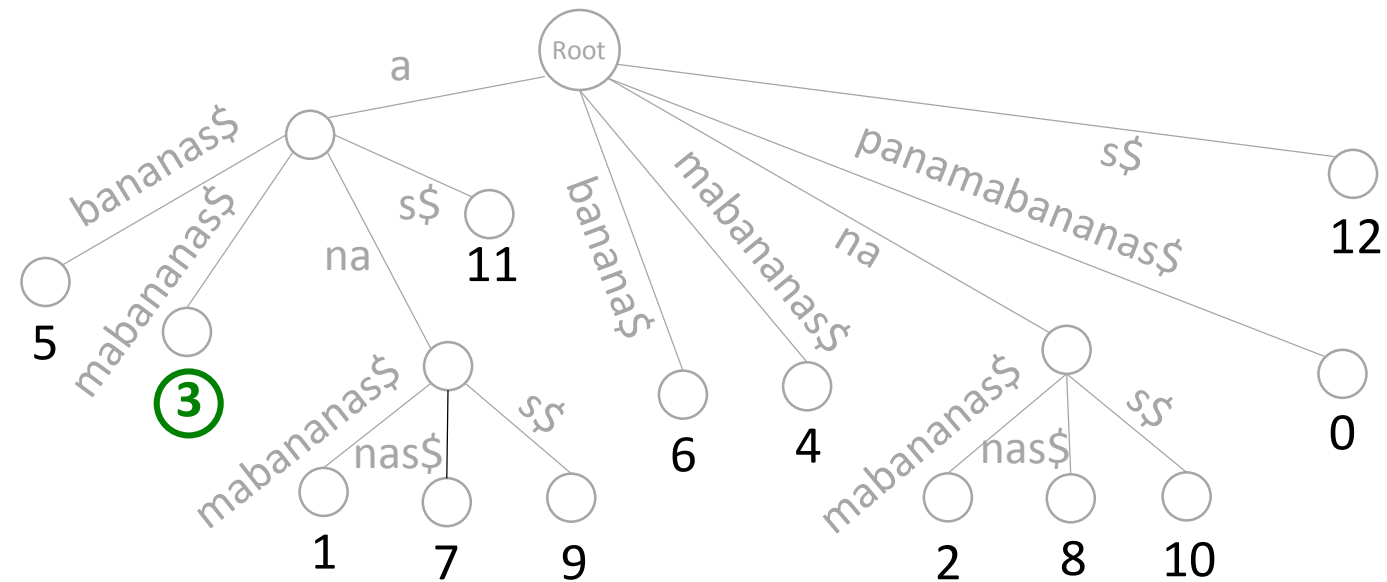
# From Suffix Tree to Suffix Array: Depth-First Traversal

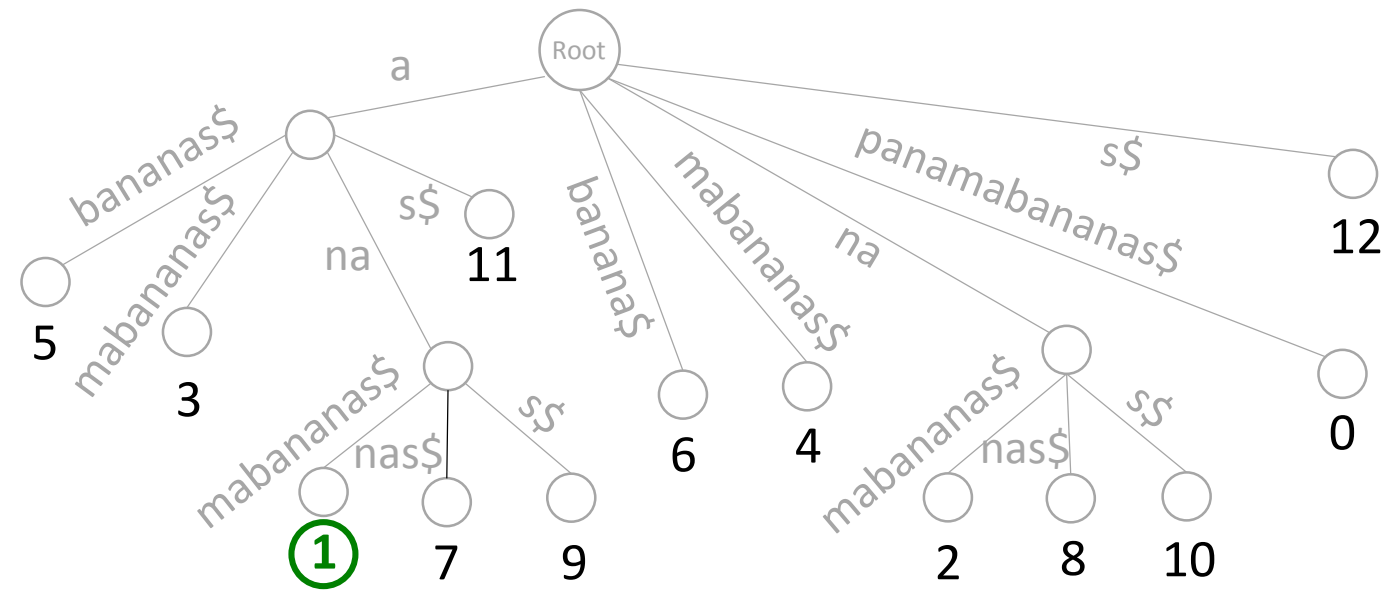

[13  5  ③  1  7  9  11  6  4  2  8  10  0  12]

# From Suffix Tree to Suffix Array: Depth-First Traversal



[13  5  3  (1)  7  9  11  6  4  2  8  10  0  12]

# From Suffix Tree to Suffix Array



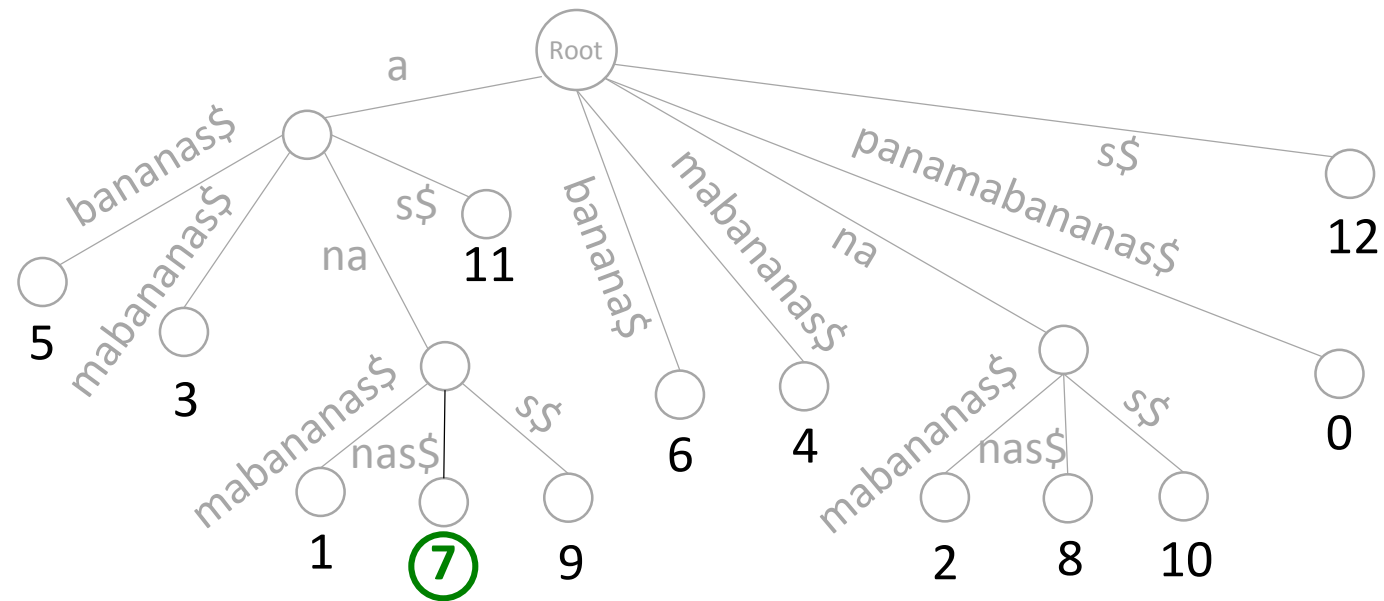[13  5  3  1  (7)  9  11  6  4  2  8  10  0  12]

# From Suffix Tree to Suffix Array



[13  5  3  1  7  (9)  11  6  4  2  8  10  0  12]

# From Suffix Tree to Suffix Array



[13  5  3  1  7  9  (11)  6  4  2  8  10  0  12]

# From Suffix Tree to Suffix Array



[13  5  3  1  7  9  11  (6)  4  2  8  10  0  12]

# From Suffix Tree to Suffix Array



[13  5  3  1  7  9  11  6  (4)  2  8  10  0  12]

# From Suffix Tree to Suffix Array



[13  5  3  1  7  9  11  6  4  ②  8  10  0  12]

# From Suffix Tree to Suffix Array
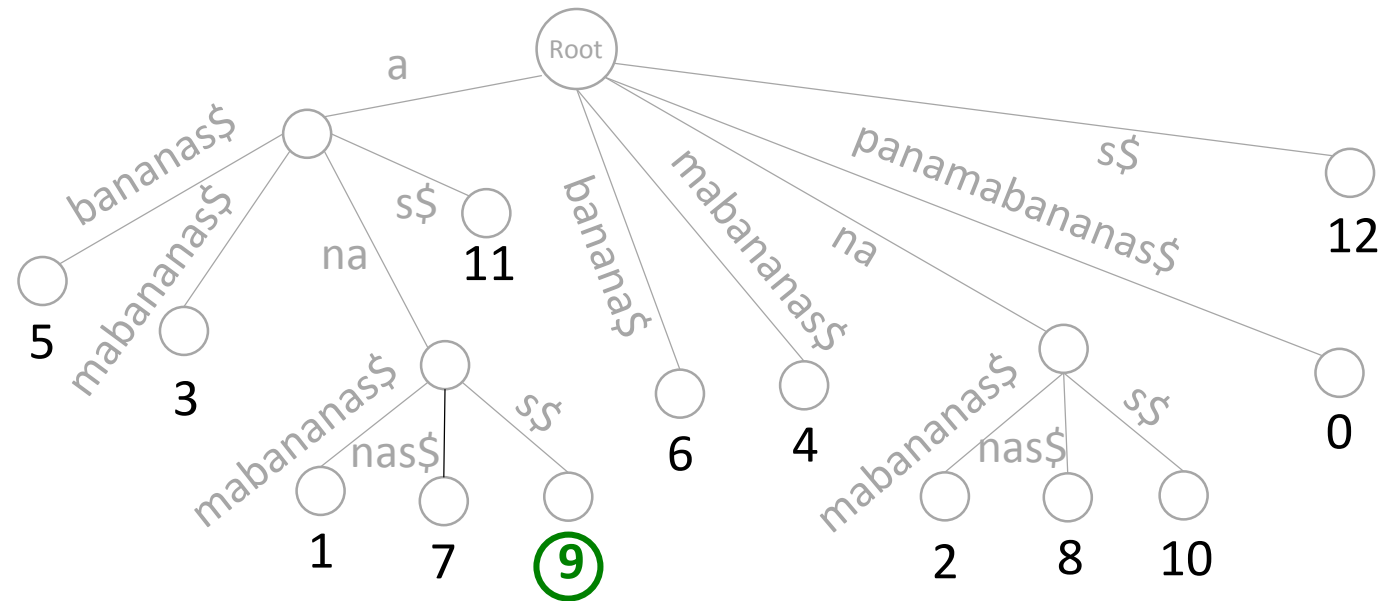


[13  5  3  1  7  9  11  6  4  2  (8)  10  0  12]

# From Suffix Tree to Suffix Array

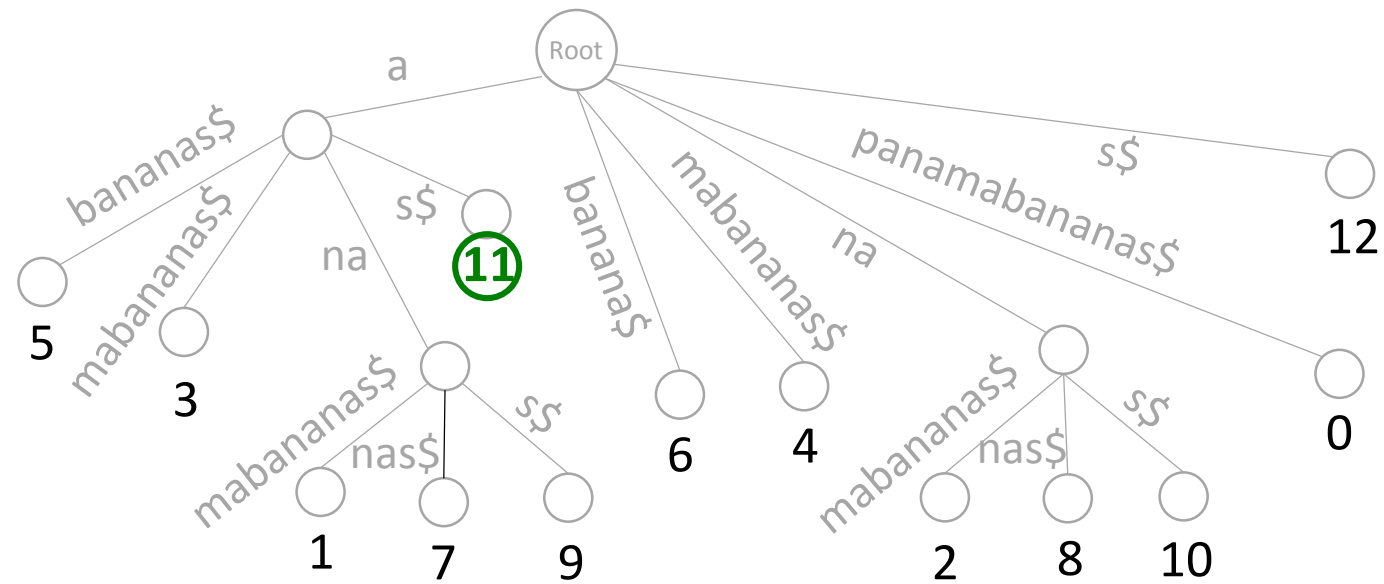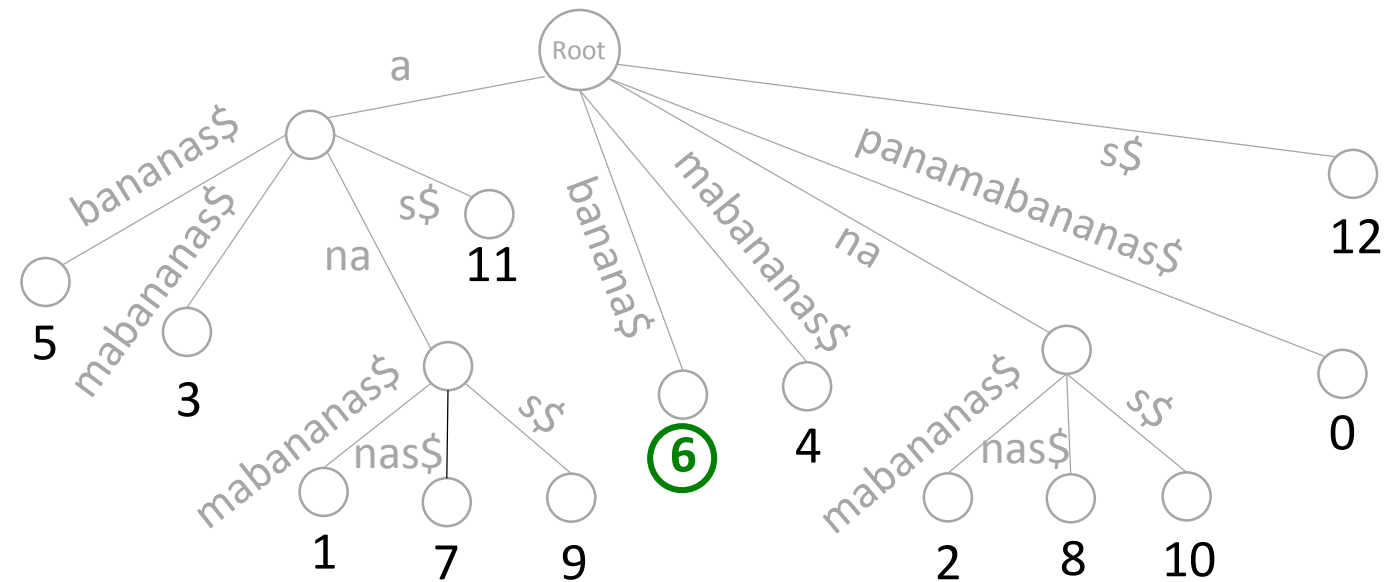

[13  5  3  1  7  9  11  6  4  2  8  (10)  0  12]

# From Suffix Tree to Suffix Array



[13    5    3    1    7    9    11    6    4    2    8    10    ⓪    12]
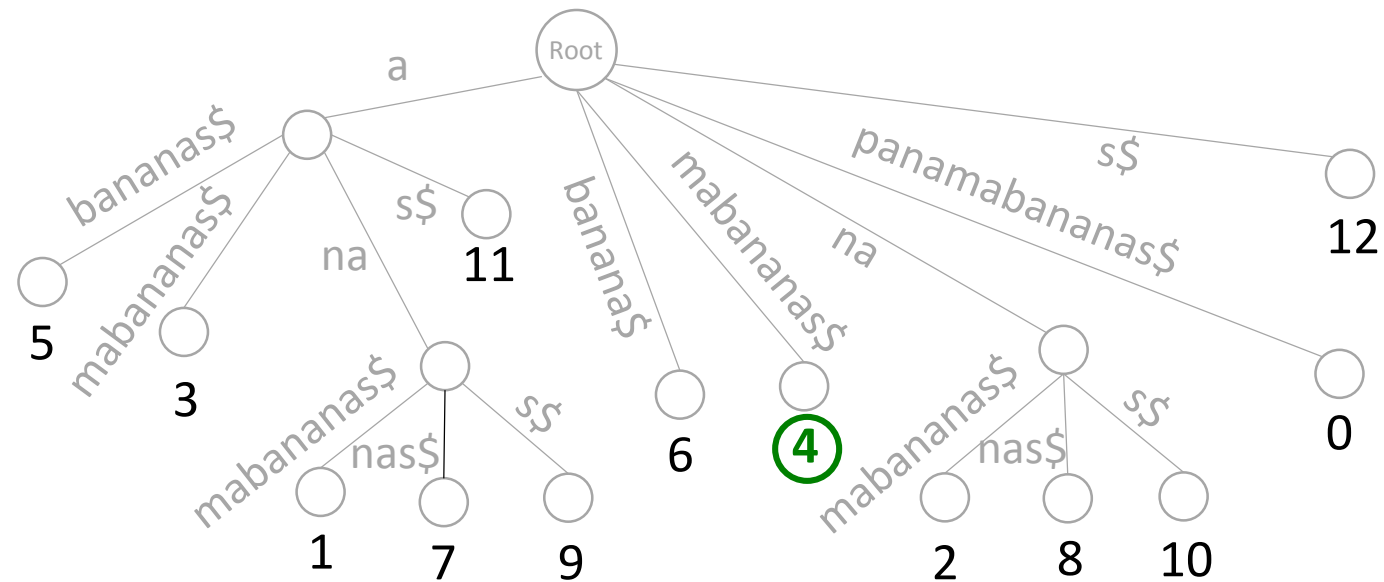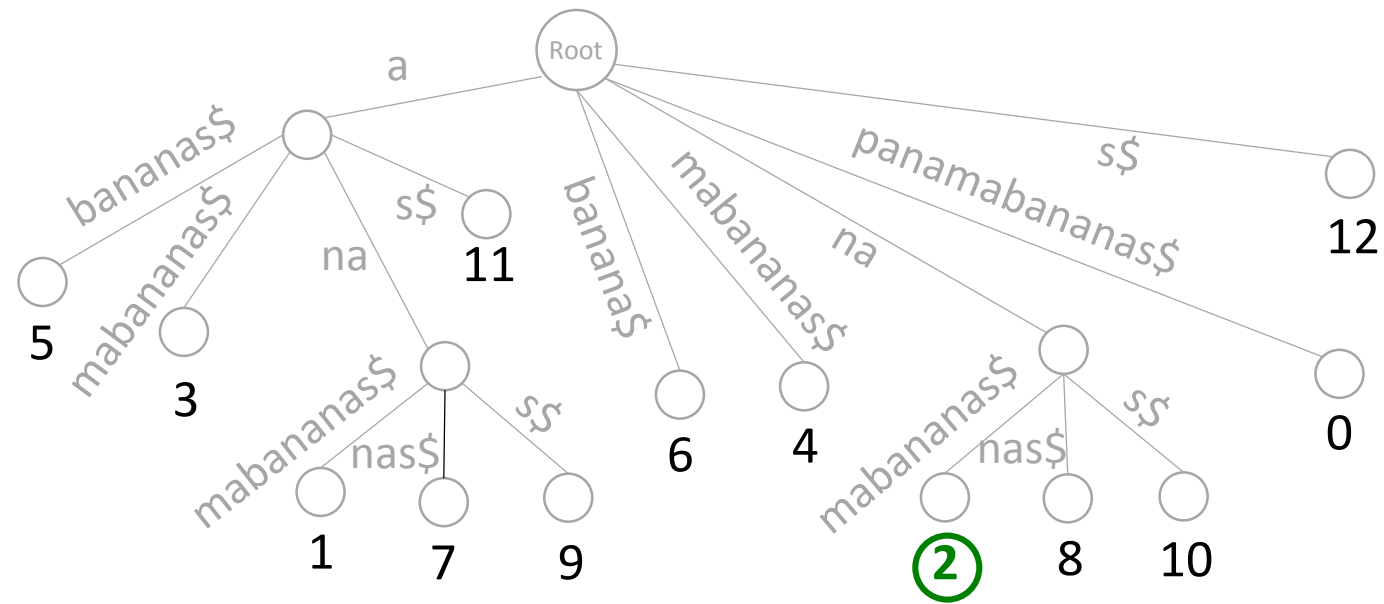
# From Suffix Tree to Suffix Array



[13   5   3   1   7   9   11   6   4   2   8   10   0   (12)]

# Constructing Suffix Array

- Depth-first traversal of suffix tree
  - O($|Text|$) time and ~20•$|Text|$ space

- Manber-Myers algorithm (1990):
  - O($|Text|$) time and ~4•$|Text|$ space

- But memory footprint is still large for human genome!

We will learn how to quickly construct suffix array
without relying on suffix tree later in this course

# Reducing Memory Footprint for Suffix Array

- Can we store only a fraction of the suffix array but still do fast pattern matching?

```
1 3
  5
  3
  1
  7
  9
1 1
  6
  4
  2
  8
1 0
  0
1 2
```

# Reducing Memory Footprint for Suffix Array

- Can we store only a fraction of the suffix array but still do fast pattern matching?

- Partial suffix array SuffixArray$_K$(*Text*) only contains values that are multiples of some integer *K*

5

1 0
0

# Using the Suffix Array to Find Matches

$\$_1$panamabananas$_1$    13

$a_1$bananas\$panam$_1$    5

$a_2$mabananas\$pan$_1$    3

**$a_3$na**mabananas\$p$_1$    **1**

**$a_4$na**nas\$panamab$_1$    **7**

**$a_5$na**s\$panamaban$_2$    **9**

$a_6$s\$panamabanan$_3$    11

$b_1$ananas\$panama$_1$    6

$m_1$abananas\$pana$_2$    4

$n_1$amabananas\$pa$_3$    2

$n_2$anas\$panamaba$_4$    8

$n_3$as\$panamabana$_5$    10

$p_1$anamabananas$\$_1$    0

$s_1$\$panamabanana$_6$    12

# Using the Partial Suffix Array to Find Matches

$\$_1$ p a n a m a b a n a n a $s_1$

$a_1$ b a n a n a s \$ p a n a $m_1$      **5**

$a_2$ m a b a n a n a s \$ p a $n_1$

**$a_3$ n a** m a b a n a n a s \$ p $_1$

**$a_4$ n a** n a s \$ p a n a m a b $_1$

**$a_5$ n a** s \$ p a n a m a b a $n_2$

$a_6$ s \$ p a n a m a b a n a $n_3$

$b_1$ a n a n a s \$ p a n a m $a_1$

$m_1$ a b a n a n a s \$ p a n $a_2$

$n_1$ a m a b a n a n a s \$ p $a_3$

$n_2$ a n a s \$ p a n a m a b $a_4$

$n_3$ a s \$ p a n a m a b a n $a_5$      **10**

$p_1$ a n a m a b a n a n a s \$ $_1$      **0**

$s_1$ \$ p a n a m a b a n a n $a_6$

# Using the Partial Suffix Array to Find Matches

$\$_1$ p a n a m a b a n a n a s $_1$

a $_1$ b a n a n a s $ p a n a m $_1$    **5**

a $_2$ m a b a n a n a s $ p a n $_1$

**a $_3$ n a** m a b a n a n a s $ p $_1$

**a $_4$ n a** n a s $ p a n a m a b $_1$    Where are these **ana** prefixes located in *Text*???

**a $_5$ n a** s $ p a n a m a b a n $_2$

a $_6$ s $ p a n a m a b a n a n $_3$

b $_1$ a n a n a s $ p a n a m a $_1$

m $_1$ a b a n a n a s $ p a n a $_2$

n $_1$ a m a b a n a n a s $ p a $_3$

n $_2$ a n a s $ p a n a m a b a $_4$

n $_3$ a s $ p a n a m a b a n a $_5$    **1 0**

p $_1$ a n a m a b a n a n a s $ $_1$    **0**

s $_1$ $ p a n a m a b a n a n a $_6$

# Focus on $a_4$na

$\$_1$panamabananas$_1$

$a_1$bananas\$panam$_1$ — 5

$a_2$mabananas\$pan$_1$

**$a_3$**namabananas\$p$_1$

**$a_4$na**nas\$panama**$b_1$** — Where is **$a_4$na**?

**$a_5$na**s\$panaman$_2$

$a_6$s\$panamanan$_3$

$b_1$ananas\$panama$_1$

$m_1$abananas\$pana$_2$

$n_1$amabananas\$pa$_3$

$n_2$anas\$panamaba$_4$

$n_3$as\$panamaban$_5$ — 10

$p_1$anamabananas$\$_1$ — 0

$s_1$\$panamabanana$_6$

# Focus on $b_1$ana

```
$1panamabananas1
a1bananas$panam1                    5
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamabanan3
b1ananas$panama1    Where is b1ana?
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamaban a5              10
p1anamabananas$1               0
s1$panamabanana6
```

# Focus on $a_1$bana

$\$_1$panamabananas$_1$

$a_1b$ananas$pana$m_1$     Where is $a_1$**bana**?     **5**

$a_2$mabananas$pa$n_1$

$a_3$namabananas$p$_1$

$a_4$**na**nas$panama$b_1$

$a_5$**na**s$panamaba$n_2$

$a_6$s$panamaban$a$n_3$

$b_1$**ana**nas$panam$a_1$

$m_1$abananas$pana$a_2$

$n_1$amabananas$pa$a_3$

$n_2$anas$panamab$a_4$

$n_3$as$panamaban$a_5$     **10**

$p_1$anamabananas$\$_1$     **0**

$s_1$$panamabanan$a_6$

# Partial suffix array reveals position of $a_1$**bana**



partial
suffix
array

$a_1$**bana** is at position 5          **5**

$a_4$**na** is at position 7          7

$b_1$**ana** is at position 6          6

10
0

# Outline

- [Burrows-Wheeler Transform](#)

- [Inverting Burrows-Wheeler Transform](#)

- [Using BWT for Pattern Matching](#)

- [Suffix Arrays](#)

- **[Approximate Pattern Matching](#)**

# Returning to Search for Mutations

- **Approximate Pattern Matching Problem**:
  - **Input**: A string *Pattern*, a string *Text*, and an integer *d*.
  - **Output:** All positions in *Text* where the string *Pattern* appears as a substring with at most *d* mismatches.

# Revealing Mutations by Analyzing Billions of Reads

- **Multiple Approximate Pattern Matching Problem**
  - **Input**: A **set** of strings *Patterns*, a string *Text*, and an integer *d*.
  - **Output:** All positions in *Text* where a string from *Patterns* appears as a substring with at most *d* mismatches.

# BWT Saves the Day Again

- **searching for** `ana` **in** `panamabananas`

$_1$panamabanana$s_1$
$a_1$bananas\$pana$m_1$
$a_2$mabananas\$pa$n_1$
$a_3$namabananas\$$p_1$
$a_4$nanas\$panama$b_1$
$a_5$nas\$panamaba$n_2$
$a_6$s\$panamabana$n_3$
$b_1$ananas\$panam$a_1$
$m_1$abananas\$pan$a_2$
$n_1$amabananas\$p$a_3$
$n_2$anas\$panamaba$a_4$
$n_3$as\$panamabana$a_5$
$p_1$anamabananas$\$_1$
$s_1$\$panamabanan$a_6$

# BWT Saves the Day Again

- searching for an**a** in panamabananas

$\$_1$panamabanana$s_1$
**a**$_1$bananas\$pana$m_1$
**a**$_2$mabananas\$pa$n_1$
**a**$_3$namabananas\$$p_1$
**a**$_4$nanas\$panama$b_1$
**a**$_5$nas\$panamaba$n_2$
**a**$_6$s\$panamabana$n_3$
$b_1$ananas\$panam$a_1$
$m_1$abananas\$pan$a_2$
$n_1$amabananas\$p$a_3$
$n_2$anas\$panamab$a_4$
$n_3$as\$panamabana$a_5$
$p_1$anamabananas$\$_1$
$s_1$\$panamabanan$a_6$

# BWT Saves the Day Again

- searching for a**na** in panamabananas

$\$_1$ p a n a m a b a n a n a s $_1$
**a**$_1$ b a n a n a s \$ p a n a **m**$_1$
**a**$_2$ m a b a n a n a s \$ p a **n**$_1$
**a**$_3$ n a m a b a n a n a s \$ **p**$_1$
**a**$_4$ n a n a s \$ p a n a m a **b**$_1$
**a**$_5$ n a s \$ p a n a m a b a **n**$_2$
**a**$_6$ s \$ p a n a m a b a n a **n**$_3$
b$_1$ a n a n a s \$ p a n a m a$_1$
m$_1$ a b a n a n a s \$ p a n a$_2$
n$_1$ a m a b a n a n a s \$ p a$_3$
n$_2$ a n a s \$ p a n a m a b a$_4$
n$_3$ a s \$ p a n a m a b a n a$_5$
p$_1$ a n a m a b a n a n a s \$$_1$
s$_1$ \$ p a n a m a b a n a n a$_6$

Exact matching

# BWT Pattern Matching with 1 Mismatch

- searching for a**na** in panamabananas

To allow for 1 mismatch, we need to analyze the rows ending in red letters as well.

$\$_1$panamabanan$s_1$
$\mathbf{a}_1$bananas\$pana$\mathbf{m}_1$
$\mathbf{a}_2$mabananas\$pan$\mathbf{n}_1$
$\mathbf{a}_3$namabananas\$$\mathbf{p}_1$
$\mathbf{a}_4$nanas\$panama$\mathbf{b}_1$
$\mathbf{a}_5$nas\$panamaba$\mathbf{n}_2$
$\mathbf{a}_6$s\$panamaban$\mathbf{n}_3$
$b_1$ananas\$panam$a_1$
$m_1$abananas\$pan$a_2$
$n_1$amabananas\$p$a_3$
$n_2$anas\$panamab$a_4$
$n_3$as\$panamaban$a_5$
$p_1$anamabananas$\$_1$
$s_1$\$panamabanan$a_6$

Approximate matching with at most 1 mismatch

# BWT Pattern Matching with 1 Mismatch

- searching for a**na** in panamabananas

To allow for 1 mismatch, we need to analyze the rows ending in red letters as well.

# Mismatches

```
$₁panamabananas₁
a₁bananas$panam₁      1
a₂mabananas$pan₁      0
a₃namabananas$p₁      1
a₄nanas$panamab₁      1
a₅nas$panamaban₂      0
a₆s$panamabanan₃      0
b₁ananas$panama₁
m₁abananas$pana₂
n₁amabananas$pa₃
n₂anas$panamaba₄
n₃as$panamaban a₅
p₁anamabananas$₁
s₁$panamabanan a₆
```

# BWT Pattern Matching with 1 Mismatch

- searching for a**na** in panamabananas

Now we analyze all rows with at most 1 mismatch using the First-Last property.

# Mismatches

$$\$_1 \, p \, a \, n \, a \, m \, a \, b \, a \, n \, a \, n \, a \, s_1$$

| | # Mismatches |
|---|---|
| $\textbf{a}_1 \, b \, a \, n \, a \, n \, a \, s \, \$ \, p \, a \, n \, a \, \textbf{m}_1$ | 1 |
| $\textbf{a}_2 \, m \, a \, b \, a \, n \, a \, n \, a \, s \, \$ \, p \, a \, \textbf{n}_1$ | 0 |
| $\textbf{a}_3 \, n \, a \, m \, a \, b \, a \, n \, a \, n \, a \, s \, \$ \, \textbf{p}_1$ | 1 |
| $\textbf{a}_4 \, n \, a \, n \, a \, s \, \$ \, p \, a \, n \, a \, m \, a \, \textbf{b}_1$ | 1 |
| $\textbf{a}_5 \, n \, a \, s \, \$ \, p \, a \, n \, a \, m \, a \, b \, a \, \textbf{n}_2$ | 0 |
| $\textbf{a}_6 \, s \, \$ \, p \, a \, n \, a \, m \, a \, b \, a \, n \, a \, \textbf{n}_3$ | 0 |

$$b_1 \, a \, n \, a \, n \, a \, s \, \$ \, p \, a \, n \, a \, m \, a_1$$
$$m_1 \, a \, b \, a \, n \, a \, n \, a \, s \, \$ \, p \, a \, n \, a_2$$
$$n_1 \, a \, m \, a \, b \, a \, n \, a \, n \, a \, s \, \$ \, p \, a_3$$
$$n_2 \, a \, n \, a \, s \, \$ \, p \, a \, n \, a \, m \, a \, b \, a_4$$
$$n_3 \, a \, s \, \$ \, p \, a \, n \, a \, m \, a \, b \, a \, n \, a_5$$
$$p_1 \, a \, n \, a \, m \, a \, b \, a \, n \, a \, n \, a \, s \, \$_1$$
$$s_1 \, \$ \, p \, a \, n \, a \, m \, a \, b \, a \, n \, a \, n \, a_6$$

# BWT Pattern Matching with 1 Mismatch

- searching for a**na** in panamabananas

Now we analyze all rows with at most 1 mismatch using the First-Last property.

# Mismatches

| | |
|---|---|
| $\$_1$panamabanana$s_1$ | |
| $a_1$bananas\$pana$\textbf{m}_1$ | 1 |
| $a_2$mabananas\$pa$\textbf{n}_1$ | 0 |
| $a_3$namabananas\$$\textbf{p}_1$ | 1 |
| $a_4$nanas\$panama$\textbf{b}_1$ | 1 |
| $a_5$nas\$panamaba$\textbf{n}_2$ | 0 |
| $a_6$s\$panamabana$\textbf{n}_3$ | 0 |
| $\textbf{b}_1\textbf{a}$nanas\$panam$a_1$ | |
| $\textbf{m}_1\textbf{a}$bananas\$pan$a_2$ | |
| $\textbf{n}_1\textbf{a}$mabananas\$p$a_3$ | |
| $\textbf{n}_2\textbf{a}$nas\$panamab$a_4$ | |
| $\textbf{n}_3\textbf{a}$s\$panamaban$a_5$ | |
| $\textbf{p}_1\textbf{a}$namabananas$\$_1$ | |
| $s_1$\$panamabanan$a_6$ | |

# BWT Pattern Matching with 1 Mismatch

- searching for a**na** in panamabananas

Now we analyze all rows with at most 1 mismatch using the First-Last property.

| | |
|---|---|
| $\$_1$ p a n a m a b a n a n a s $_1$ | |
| a $_1$ b a n a n a s $ p a n a m $_1$ | |
| a $_2$ m a b a n a n a s $ p a n $_1$ | |
| a $_3$ n a m a b a n a n a s $ p $_1$ | |
| a $_4$ n a n a s $ p a n a m a b $_1$ | |
| a $_5$ n a s $ p a n a m a b a n $_2$ | |
| a $_6$ s $ p a n a m a b a n a n $_3$ | |
| **b $_1$ a** n a n a s $ p a n a m a $_1$ | 1 |
| **m $_1$ a** b a n a n a s $ p a n a $_2$ | 1 |
| **n $_1$ a** m a b a n a n a s $ p a $_3$ | 0 |
| **n $_2$ a** n a s $ p a n a m a b a $_4$ | 0 |
| **n $_3$ a** s $ p a n a m a b a n a $_5$ | 0 |
| **p $_1$ a** n a m a b a n a n a s $_1$ | 1 |
| s $_1$ $ p a n a m a b a n a n a $_6$ | |

# BWT Pattern Matching with 1 Mismatch

- searching for **ana** in `panamabananas`

# Mismatches

$\$_1$`panamabanas`$_1$

$a_1$`bananas$pana`$m_1$

$a_2$`mabananas$pa`$n_1$

$a_3$`namabananas$`$p_1$

$a_4$`nanas$panama`$b_1$

$a_5$`nas$panamaban`$n_2$

$a_6$`s$panamabana`$n_3$

**$b_1$** **a**`nanas$panam`**$a_1$**  1

**$m_1$** **a**`bananas$pan`**$a_2$**  1

**$n_1$** **a**`mabananas$p`**$a_3$**  0

**$n_2$** **a**`nas$panamab`**$a_4$**  0

**$n_3$** **a**`s$panamaban`**$a_5$**  0

This row results in a 2nd mismatch (the $), so we discard it.

**$p_1$** **a**`namabananas`**$\$_1$**  **2**

$s_1$`$panamabanan`$a_6$

# Five Approximate Matches Found!

- searching for **ana** in panamabananas

# Mismatches

$\$_1$panamanas$s_1$

$a_1ba$nanas$\$$pana$m_1$     1

$a_2ma$bananas$\$$pan$n_1$     1

$a_3na$mabananas$\$$p$p_1$     0

$a_4na$nas$\$$panama$b_1$     0

$a_5na$s$\$$panamaban$n_2$     0

$a_6$s$\$$panamabanan$n_3$

$b_1$ananas$\$$panam$a_1$

$m_1$abananas$\$$pan$a_2$

$n_1$amabananas$\$$p$a_3$

$n_2$anas$\$$panamab$a_4$

$n_3$as$\$$panamaban$a_5$

$p_1$anamabananas$\$_1$

$s_1\$$panamabanan$a_6$

# Where Are The Matches?

- searching for **ana** in panamabananas

Suffix Array

$\$_1$ p a n a m a b a n a n a s $_1$

$a_1$ $b$ $a$ n a n a s $ p a n a m $_1$       5

$a_2$ $m$ $a$ b a n a n a s $ p a n $_1$       3

$a_3$ $n$ $a$ m a b a n a n a s $ p $_1$       1

$a_4$ $n$ $a$ n a s $ p a n a m a b $_1$       7

$a_5$ $n$ $a$ s $ p a n a m a b a n $_2$       9

$a_6$ s $ p a n a m a b a n a n $_3$

b $_1$ a n a n a s $ p a n a m a $_1$

m $_1$ a b a n a n a s $ p a n a $_2$

n $_1$ a m a b a n a n a s $ p a $_3$

n $_2$ a n a s $ p a n a m a b a $_4$

n $_3$ a s $ p a n a m a b a n a $_5$

p $_1$ a n a m a b a n a n a s $ $_1$

s $_1$ $ p a n a m a b a n a n a $_6$

In reality, approximate pattern matching with BWT is more complex (we omitted various details)