# Suffix Array

A suffix array is a data structure designed for efficient searching of a large text. It is an array of integers giving the starting position of the suffixes of a string in lexicographical order.

It can be used as index to quickly locate every occurrence of a substring within the string. Finding every occurrence of the substring is equivalent to finding every suffix that begins with the substring.

Thanks to the lexicographical ordering these suffixes will be grouped together in the suffix array, and can be found efficiently using a binary search.

# Suffix Array

| Text | a | b | r | a | c | a | d | a | b | r | a |
|------|---|---|---|---|---|---|---|---|---|---|---|
| Index | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |

**Enumerate all the suffixes**

| Suffix | Index |
|--------|-------|
| a b r a c a d a b r a | 0 |
| b r a c a d a b r a | 1 |
| r a c a d a b r a | 2 |
| a c a d a b r a | 3 |
| c a d a b r a | 4 |
| a d a b r a | 5 |
| d a b r a | 6 |
| a b r a | 7 |
| b r a | 8 |
| r a | 9 |
| a | 10 |

**Sort the suffixes**

| Sorted Suffix | Index |
|---------------|-------|
| a | 10 |
| a b r a | 7 |
| a b r a c a d a b r a | 0 |
| a c a d a b r a | 3 |
| a d a b r a | 5 |
| b r a | 8 |
| b r a c a d a b r a | 1 |
| c a d a b r a | 4 |
| d a b r a | 6 |
| r a | 9 |
| r a c a d a b r a | 2 |

**Output**

| 10 | 7 | 0 | 3 | 5 | 8 | 1 | 4 | 6 | 9 | 2 |
|----|---|---|---|---|---|---|---|---|---|---|

## Binary search of suffix "ra"

| Sorted Suffix | Index | |
|---------------|-------|---|
| a | 10 | |
| a b r a | 7 | |
| a b r a c a d a b r a | 0 | |
| a c a d a b r a | 3 | |
| a d a b r a | 5 | |
| b r a | 8 | ←— 1 |
| b r a c a d a b r a | 1 | |
| c a d a b r a | 4 | |
| d a b r a | 6 | ←— 2 |
| r a | 9 | ←3 |
| r a c a d a b r a | 2 | |

# Suffix array is a memory-efficient alternative to suffix trees

$$\text{SUFFIXARRAY (``panamabananas\$'')} = [13,5,3,1,7,9,11,6,4,2,8,10,0,12]$$

| Sorted Suffixes | Starting Positions |
|---|---|
| $ | 13 |
| abananas$ | 5 |
| amabananas$ | 3 |
| anamabananas$ | 1 |
| ananas$ | 7 |
| anas$ | 9 |
| as$ | 11 |
| bananas$ | 6 |
| mabananas$ | 4 |
| namabananas$ | 2 |
| nanas$ | 8 |
| nas$ | 10 |
| panamabananas$ | 0 |
| s$ | 12 |

```
BINARYSEARCH(ARRAY, key, minIndex, maxIndex)
    while maxIndex ≥ minIndex
        midIndex ← (minIndex + maxIndex)/2
        if ARRAY(midIndex) = key
            return midIndex
        else if ARRAY(midIndex) < key
            minIndex ← midIndex + 1
        else
            maxIndex ← midIndex − 1
    return "key not found"
```

# Suffix Array

**How do we build a suffix array?**

- Easiest solution: build a suffix tree

- Traverse the tree in DFS, lexicographically picking edges outgoing from each node and fill the suffix array

- O(n) time

- Waste of space: can we do it directly in O(n) time? Unknown until 2003

# Idea #1: Run-Length Encoding

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

*Genome*

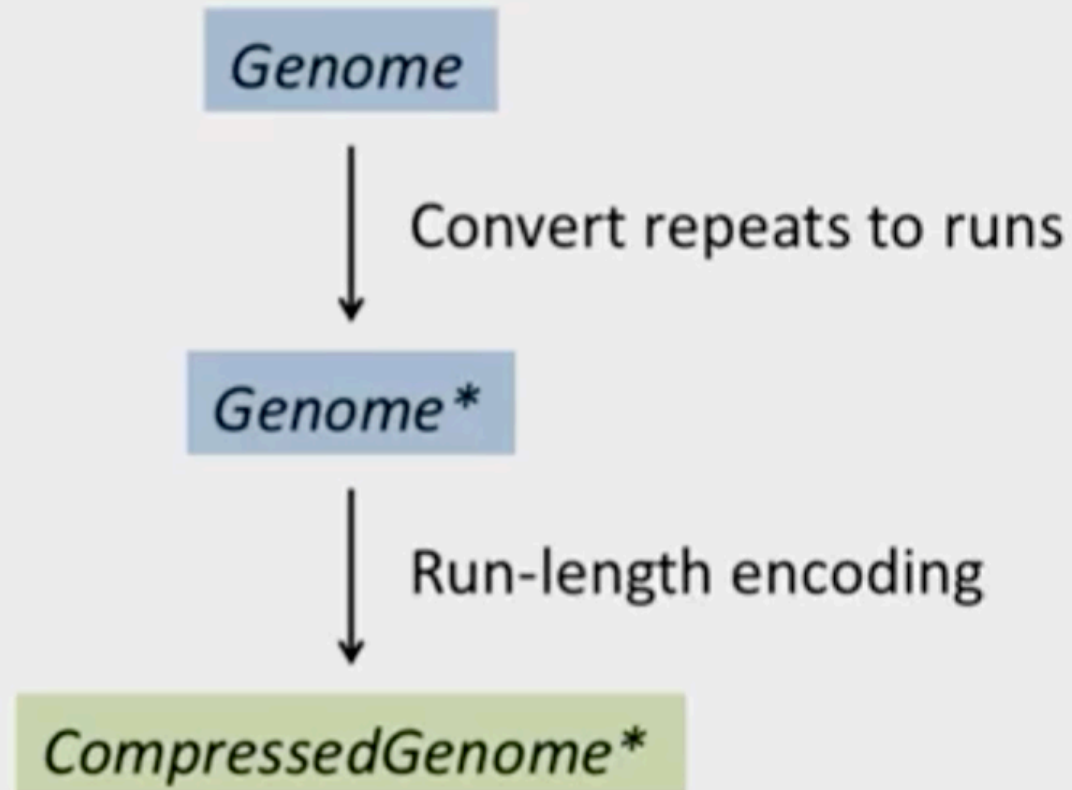GGGGGGGGGGCCCCCCCCCCCAAAAAAATTTTTTTTTTTTTTTCCCCCG

↓

10G11C7A15T5C1G

Run-length encoding

# Converting Repeats to Runs

- ...but they do have lots of repeats!

*Genome*

$\downarrow$ Convert repeats to runs

*Genome\**
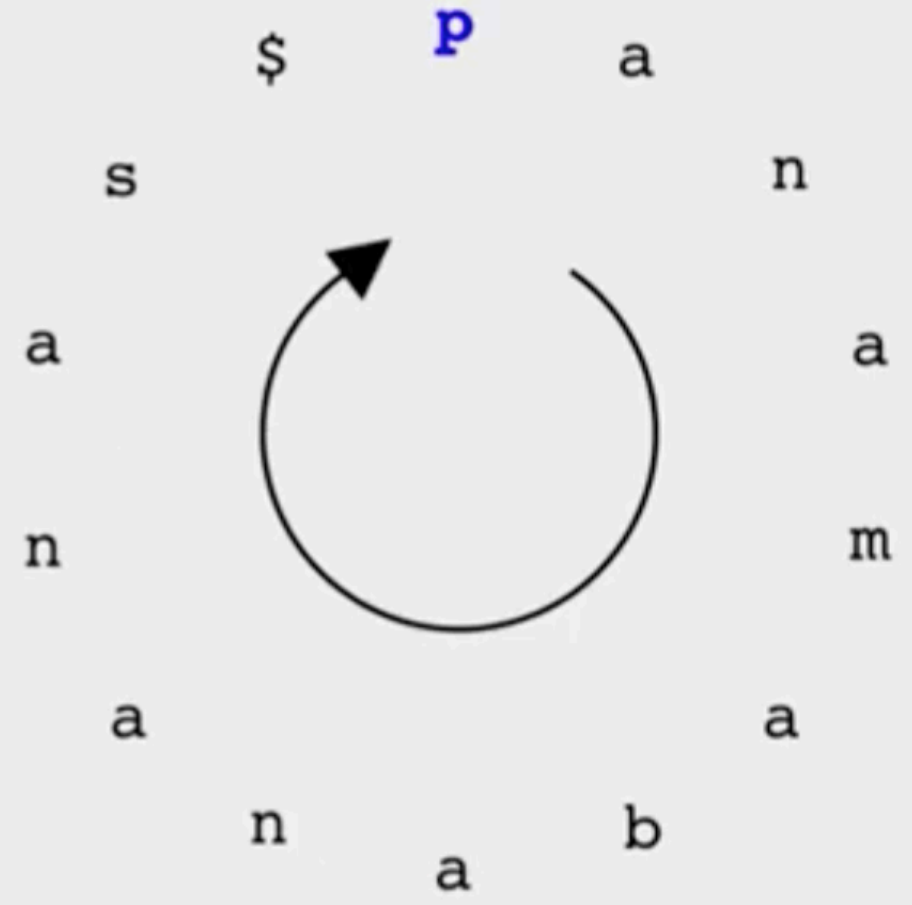
$\downarrow$ Run-length encoding

*CompressedGenome\**

# Burrows-Wheeler transformation

It performs a transformation of an input sequence consisting of a reversible permutation of the sequence characters which gives a new string that is **easier to compress**: if the original string has several substrings that occur often, then the transformed string has several places where a single character is repeated multiple time.

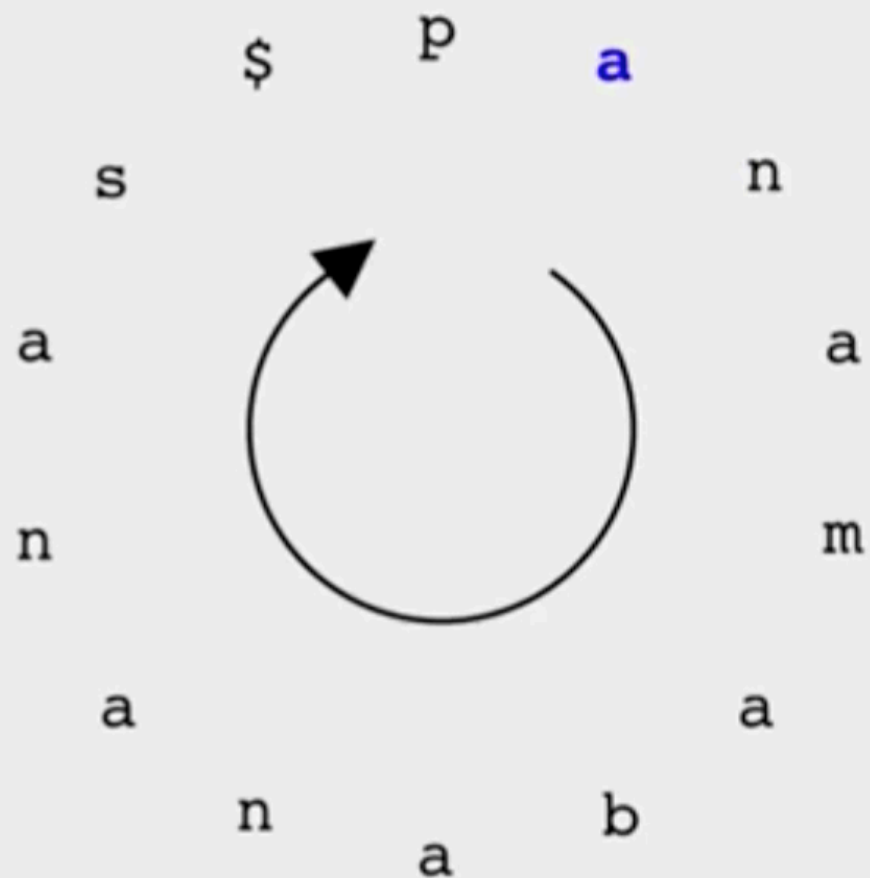# The **B**urrows-**W**heeler Transform

panamabananas$



Form all cyclic rotations of
  "panamabananas$"

# The Burrows-Wheeler Transform

```
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
```

Form all cyclic rotations of
"panamabananas$"

# The Burrows-Wheeler Transform

```
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
```
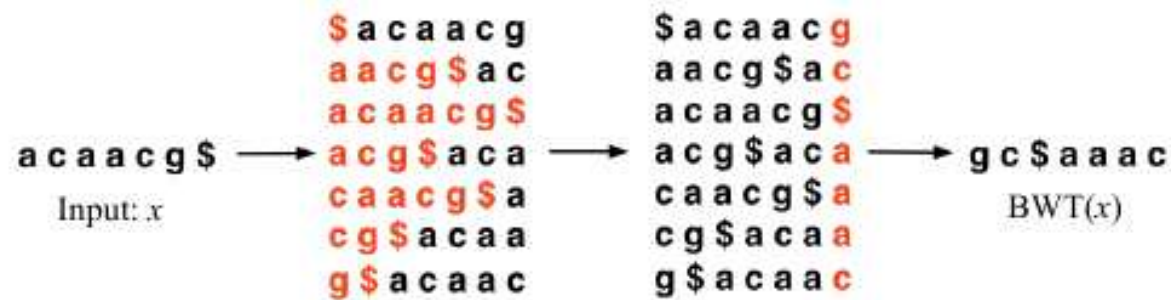
Form all cyclic rotations of          Sort the strings
"panamabananas$"          lexicographically
                                      ($ comes first)

Reversible permutation used originally in lossless data compression:



Input $(x)$ transformed in the matrix of all circular shifts of $x$, sorted lexicographically. Then the last column of the matrix became the BWT$(x)$

# Another Example

appellee$

appellee$
ppellee$a
pellee$ap
ellee$app                    $appellee
llee$appe      sort          appellee$
lee$appel      →             e$appelle
ee$appell                    ee$appell
e$appelle                    ellee$app
$appellee                    lee$appel
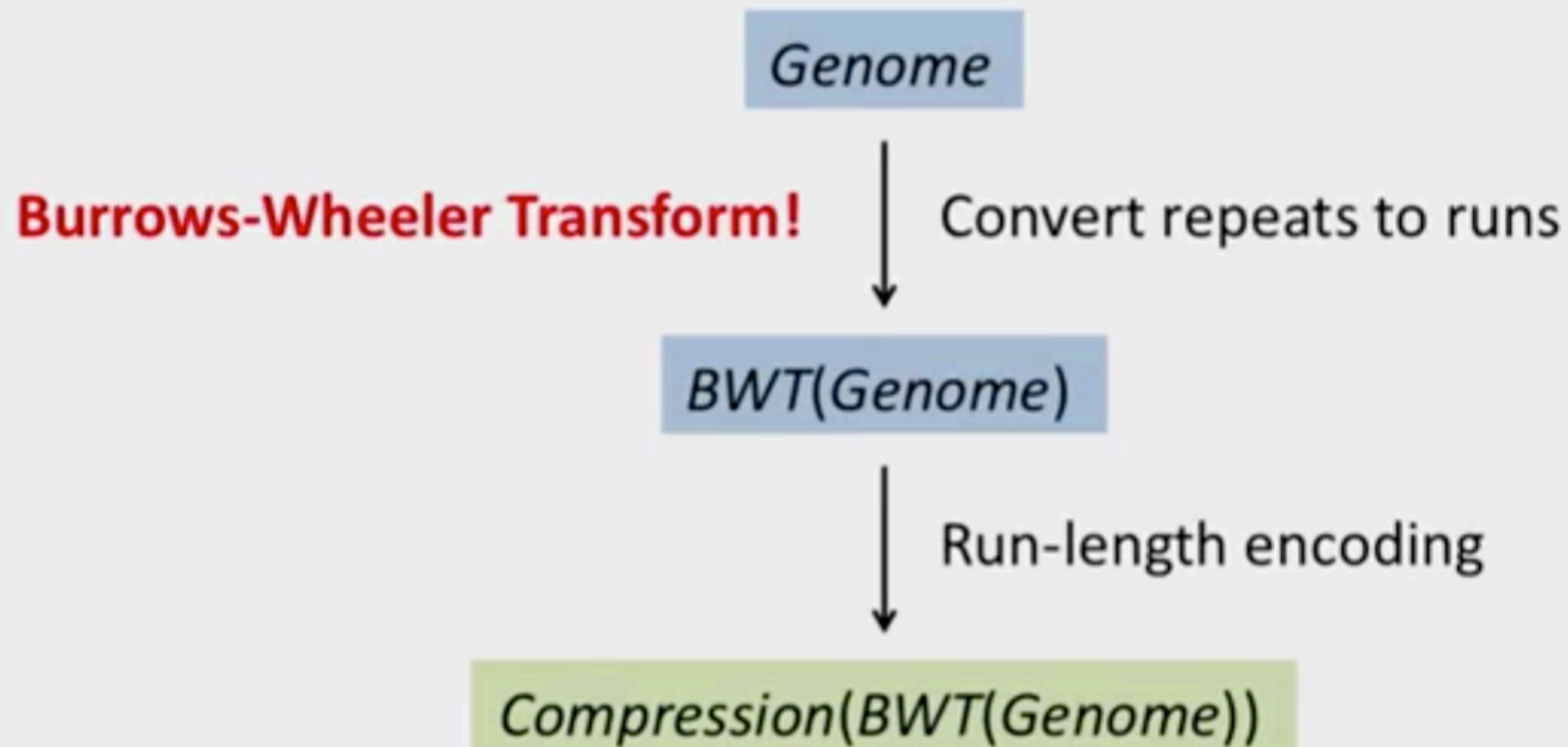                             llee$appe
                             pellee$ap
                             ppellee$a
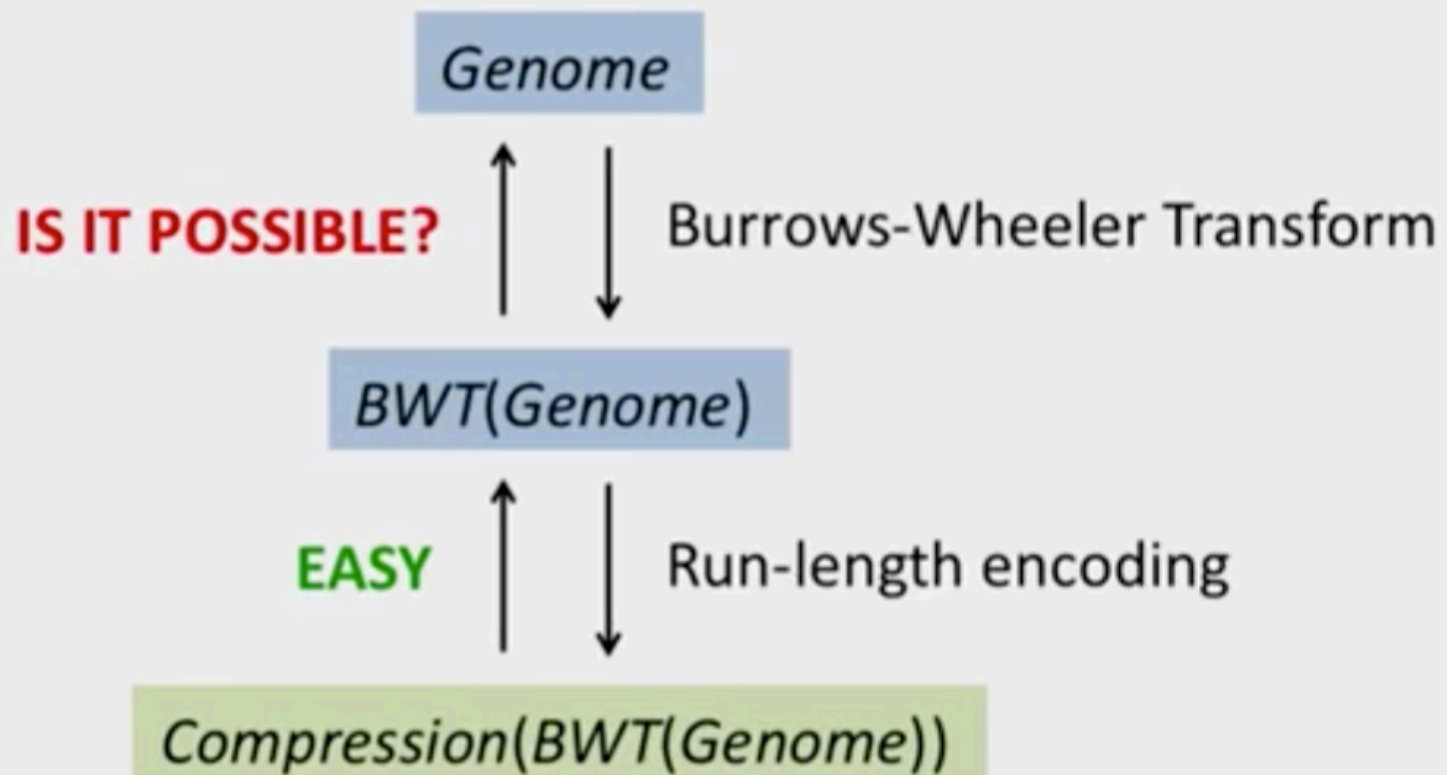
BWT(appellee$) =
e$elplepa

Doesn't always improve
the compressibility…

# BWT: Converting Repeats to Runs

Genome

**Burrows-Wheeler Transform!**      Convert repeats to runs

BWT(*Genome*)

Run-length encoding

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

# Burrows-Wheeler Transform (BWT)

- Transform: ^BANANA@ INTO: BNN^AA@A

> **function** BWT (*string* s)
>   create a table, rows are all possible rotations of s
>   sort rows alphabetically
>   **return** (last column of the table)

| All Rotations | Sorted List of Rotations | Output Last Column |
|---|---|---|
| ^BANANA@ | ANANA@^B | BNN^AA@A |
| @^BANANA | ANA@^BAN | |
| A@^BANAN | A@^BANAN | |
| NA@^BANA | BANANA@^ | |
| ANA@^BAN | NANA@^BA | |
| NANA@^BA | NA@^BANA | |
| ANANA@^B | ^BANANA@ | |
| BANANA@^ | @^BANANA | |

- Reversible

> **function** inverseBWT (*string* s)
>   create empty table
>   **repeat** length(s) times
>     insert s as a column of table before first column of the table // first insert creates first column
>     sort rows of the table alphabetically
>   **return** (row that ends with the 'EOF' character)

**Last column only** suffices to reconstruct **entire matrix**, and thus recover **original string**



| Add 1 | Sort 1 | Add 2 | Sort 2 | Add 3 | Sort 3 | Add 4 | Sort 4 | Add 5 | Sort 5 | Add 6 | Sort 6 | Add 7 | Sort 7 | Add 8 | Sort 8 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| last | 1st col | pairs | 2nd col | triples | 3rd col | 4mers | 4th col | 5mers | 5th col | 6-mers | 6th col | 7-mers | 7th col | 8-mers | Full matrix |

# Recovering the string

BWT

sort
BWT → first column

→ first 2 columns

→ first 3 columns

$appellee
appellee$
e$appelle
ee$appell
ellee$app
lee$appel
llee$appe
pellee$ap
ppellee$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 = [s₁,s₂,s₃,...,sₙ]
    for i = 1..n:
        sort B
        prepend sᵢ to B[i]
    return row of B that ends with $
```

# do$oodwg    Another BWT Example

| Prepend | Sort | Prepend | Sort | Prepend | Sort | Prepend | Sort |
|---|---|---|---|---|---|---|---|
| d $ | $d | d $d | $do | d$do | $dog | d $dog | $dogw |
| o d | d$ | o d$ | d$d | od$d | d$do | o d$do | d$dog |
| $ d | do | $ do | dog | $dog | dogw | $ dogw | dogwo |
| o g | gw | o gw | gwo | ogwo | gwoo | o gwoo | gwood |
| o o | od | o od | od$ | ood$ | od$d | o od$d | od$do |
| d o | og | d og | ogw | dogw | ogwo | d ogwo | ogwoo |
| w o | oo | w oo | ood | wood | ood$ | w ood$ | ood$d |
| g w | wo | g wo | woo | gwoo | wood | g wood | wood$ |

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

All cyclic rotation of string *panamabananas$*, note that the first row of matrix M contains the word with the special character in first position

| Cyclic Rotations | M("panamabananas$") |
|---|---|
| panamabananas$ | $ p a n a m a b a n a n a s |
| $panamabananas | a b a n a n a s $ p a n a m |
| s$panamabanana | a m a b a n a n a s $ p a n |
| as$panamabanan | a n a m a b a n a n a s $ p |
| nas$panamabana | a n a n a s $ p a n a m a b |
| anas$panamaban | a n a s $ p a n a m a b a n |
| nanas$panamaba | a s $ p a n a m a b a n a n |
| ananas$panamab | b a n a n a s $ p a n a m a |
| bananas$panama | m a b a n a n a s $ p a n a |
| abananas$panam | n a m a b a n a n a s $ p a |
| mabananas$pana | n a n a s $ p a n a m a b a |
| amabananas$pan | n a s $ p a n a m a b a n a |
| namabananas$pa | p a n a m a b a n a n a s $ |
| anamabananas$p | s $ p a n a m a b a n a n a |

If we determine the first row of M(*Text*) then we can move the $ to the end of this row to reproduce *Text*.

How do we determine the remaining symbols on this first row, if all we know is *FirstColumn* and *LastColumn*?

$ ? ? ? ? ? ? ? ? ? ? a
a ? ? ? ? ? ? ? ? ? ? r
a ? ? ? ? ? ? ? ? ? ? d
a ? ? ? ? ? ? ? ? ? ? $
a ? ? ? ? ? ? ? ? ? ? r
a ? ? ? ? ? ? ? ? ? ? c
b ? ? ? ? ? ? ? ? ? ? a
b ? ? ? ? ? ? ? ? ? ? a
c ? ? ? ? ? ? ? ? ? ? a
d ? ? ? ? ? ? ? ? ? ? a
r ? ? ? ? ? ? ? ? ? ? b
r ? ? ? ? ? ? ? ? ? ? b

Note that the first symbol in **Text** must follow $ in any cyclic rotation of **Text**.

```
$ a ? ? ? ? ? ? ? ? ? a
a ? ? ? ? ? ? ? ? ? ? r
a ? ? ? ? ? ? ? ? ? ? d
a ? ? ? ? ? ? ? ? ? ? $
a ? ? ? ? ? ? ? ? ? ? r
a ? ? ? ? ? ? ? ? ? ? c
b ? ? ? ? ? ? ? ? ? ? a
b ? ? ? ? ? ? ? ? ? ? a
c ? ? ? ? ? ? ? ? ? ? a
d ? ? ? ? ? ? ? ? ? ? a
r ? ? ? ? ? ? ? ? ? ? b
r ? ? ? ? ? ? ? ? ? ? b
```

**Rule1:** The symbol in **LastColumn** must precede the symbol of **Text** found in the same row of the **FirstColumn**

The next symbol of Text should be the first symbol in a row of M(**Text**) that end in *a*.
But, five rows and in an *a* and we don't known which of them is the correct one!

```
$ a b ? ? ? ? ? ? ? ? a      $ a c ? ? ? ? ? ? ? ? a      $ a d ? ? ? ? ? ? ? ? a
a ? ? ? ? ? ? ? ? ? ? r      a ? ? ? ? ? ? ? ? ? ? r      a ? ? ? ? ? ? ? ? ? ? r
a ? ? ? ? ? ? ? ? ? ? d      a ? ? ? ? ? ? ? ? ? ? d      a ? ? ? ? ? ? ? ? ? ? d
a ? ? ? ? ? ? ? ? ? ? $      a ? ? ? ? ? ? ? ? ? ? $      a ? ? ? ? ? ? ? ? ? ? $
a ? ? ? ? ? ? ? ? ? ? r      a ? ? ? ? ? ? ? ? ? ? r      a ? ? ? ? ? ? ? ? ? ? r
a ? ? ? ? ? ? ? ? ? ? c      a ? ? ? ? ? ? ? ? ? ? c      a ? ? ? ? ? ? ? ? ? ? c
b ? ? ? ? ? ? ? ? ? ? a      b ? ? ? ? ? ? ? ? ? ? a      b ? ? ? ? ? ? ? ? ? ? a
b ? ? ? ? ? ? ? ? ? ? a      b ? ? ? ? ? ? ? ? ? ? a      b ? ? ? ? ? ? ? ? ? ? a
c ? ? ? ? ? ? ? ? ? ? a      c ? ? ? ? ? ? ? ? ? ? a      c ? ? ? ? ? ? ? ? ? ? a
d ? ? ? ? ? ? ? ? ? ? a      d ? ? ? ? ? ? ? ? ? ? a      d ? ? ? ? ? ? ? ? ? ? a
r ? ? ? ? ? ? ? ? ? ? b      r ? ? ? ? ? ? ? ? ? ? b      r ? ? ? ? ? ? ? ? ? ? b
r ? ? ? ? ? ? ? ? ? ? b      r ? ? ? ? ? ? ? ? ? ? b      r ? ? ? ? ? ? ? ? ? ? b
```

The three possibilities for the third element of the first row of M(Text) when BWT(Text) is
*ard$rcaaaabb*. How would you choose among *b*, *c* and *d* for the second symbol of **Text**?

To determine the remaining symbols of *Text*, we need to define the **First-Last Property**. We can rank the six instances of a appear in *FirstColumn* and we say that they have ranks from 1 to 6.

```
$   p a n a m a b a n a n a s
a₁  b a n a n a s $ p a n a m
a₂  m a b a n a n a s $ p a n
a₃  n a m a b a n a n a s $ p
a₄  n a n a s $ p a n a m a b
a₅  n a s $ p a n a m a b a n
a₆  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
```

Considering `a1` in the FirstColumn

$$\texttt{a1}\textbf{bananas\$panam}$$

If we cyclically rotate this string, we obtain:

$$\textbf{panama}\texttt{a1}\textbf{bananas\$}$$

Now we can identify the positions of the other five instances:

$$\textbf{p}\texttt{a3}\textbf{n}\texttt{a2}\textbf{m}\texttt{a1}\textbf{b}\texttt{a4}\textbf{n}\texttt{a5}\textbf{n}\texttt{a6}\textbf{s\$}$$

**Where are the other five instance of a located in LastColumn?**

```
$  p a n a m a b a n a n a s        $  p a n a m a b a n a n a s
a₁ b a n a n a s $ p a n a m        a₁ b a n a n a s $ p a n a m
a₂ m a b a n a n a s $ p a n        a₂ m a b a n a n a s $ p a n
a₃ n a m a b a n a n a s $ p        a₃ n a m a b a n a n a s $ p
a₄ n a n a s $ p a n a m a b        a₄ n a n a s $ p a n a m a b
a₅ n a s $ p a n a m a b a n        a₅ n a s $ p a n a m a b a n
a₆ s $ p a n a m a b a n a n        a₆ s $ p a n a m a b a n a n
b  a n a n a s $ p a n a m a₁       b  a n a n a s $ p a n a m a₁
m  a b a n a n a s $ p a n a        m  a b a n a n a s $ p a n a₂
n  a m a b a n a n a s $ p a        n  a m a b a n a n a s $ p a₃
n  a n a s $ p a n a m a b a        n  a n a s $ p a n a m a b a₄
n  a s $ p a n a m a b a n a        n  a s $ p a n a m a b a n a₅
p  a n a m a b a n a n a s $        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        s  $ p a n a m a b a n a n a₆
```

**Rule2:** The *k-th* occurrence of symbol in **FirstColumn** and the *k-th* occurrence of symbol in **LastColumn** correspond to the same position of symbol in *Text*.

The rows are already ordered in lexicographically, so if we chop off the a from the beginning of each row the remaining strings should still be ordered lexicographically:

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $a_1$ | b | a | n | a | n | a | s | $ | p | a | n | a | m |
| $a_2$ | m | a | b | a | n | a | n | a | s | $ | p | a | n |
| $a_3$ | n | a | m | a | b | a | n | a | n | a | s | $ | p |
| $a_4$ | n | a | n | a | s | $ | p | a | n | a | m | a | b |
| $a_5$ | n | a | s | $ | p | a | n | a | m | a | b | a | n |
| $a_6$ | s | $ | p | a | n | a | m | a | b | a | n | a | n |

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

Adding a back to the end of each row should not change the lexicographic ordering of these rows:

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 $a_5$
s $ p a n a m a b a n a n $a_6$

$$p a_3 n a_2 m a_1 b a_4 n a_5 n a_6 s \$$$

The **Rule2** can be of course generalised for any possible symbol and any string Text.

# 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 1
  mabananas$pana 2
  namabananas$pa 3
  nanas$panamaba 4
  nas$panamabana 5
  panamabananas$
  s$panamabanana 6
```

These strings are sorted

Chop off **a**

Ordering
doesn't
change!

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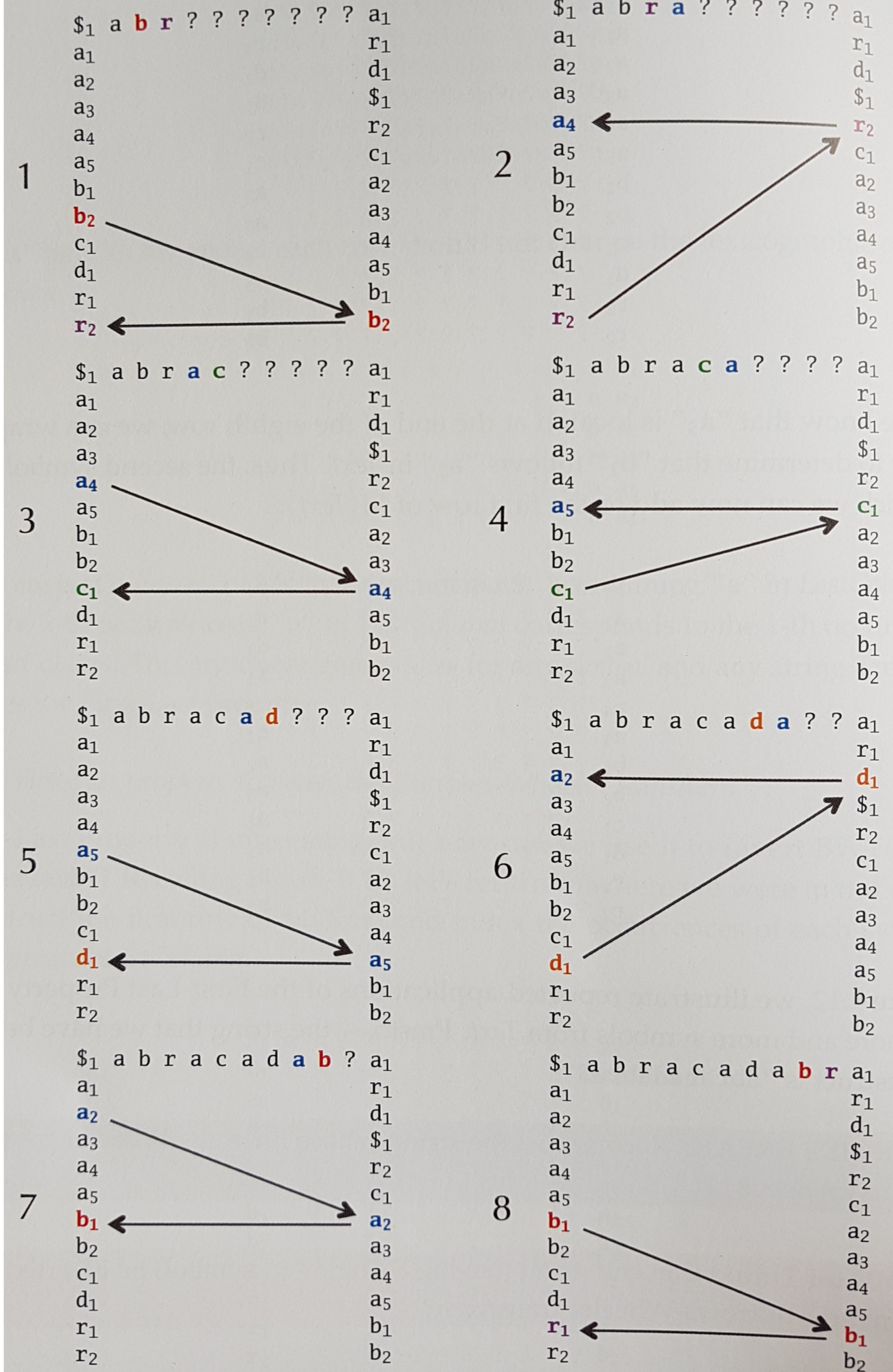
Still
sorted

Add **a**
to end

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

Still
sorted

# How do we can use the **FirstLast Property** for the **Burrows-Wheeler inversion**

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\$_1$ | a | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_1$ |
| $a_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $r_1$ |
| $a_2$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $d_1$ |
| $a_3$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $\$_1$ |
| $a_4$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $r_2$ |
| $a_5$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $c_1$ |
| $b_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_2$ |
| $b_2$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_3$ |
| $c_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_4$ |
| $d_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_5$ |
| $r_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $b_1$ |
| $r_2$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $b_2$ |

The **FL prop** reveals where **a3** is hiding in *LastColumn*:

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\$_1$ | a | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_1$ |
| $a_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $r_1$ |
| $a_2$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $d_1$ |
| $a_3$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $\$_1$ |
| $a_4$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $r_2$ |
| $a_5$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $c_1$ |
| $b_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_2$ |
| $b_2$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_3$ |
| $c_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_4$ |
| $d_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_5$ |
| $r_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $b_1$ |
| $r_2$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $b_2$ |

Since we know that **a3** is located ate the end of the eighth row, we can wrap around this row to determine that **b2** follows **a3** in the *Text*

Application of **Rule 1**

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\$_1$ | a | b | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_1$ |
| $a_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $r_1$ |
| $a_2$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $d_1$ |
| $a_3$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $\$_1$ |
| $a_4$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $r_2$ |
| $a_5$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $c_1$ |
| $b_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_2$ |
| $b_2$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_3$ |
| $c_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_4$ |
| $d_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $a_5$ |
| $r_1$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $b_1$ |
| $r_2$ | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | $b_2$ |

# Inverting the BWT

Repeated applications of the **FL prop** to reconstruct the *Text* from its **BWT**

# Pattern matching with BWT

Each row of M(*Text*) begins with a different suffix of the *Text*. Since these suffixes are already ordered lexicographically, any matches of *Pattern* in *Text* will clump together at the beginning of rows of M(*Text*).

```
$ 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
```

How do we do pattern matching without knowing the entire matrix **M(Text)?**

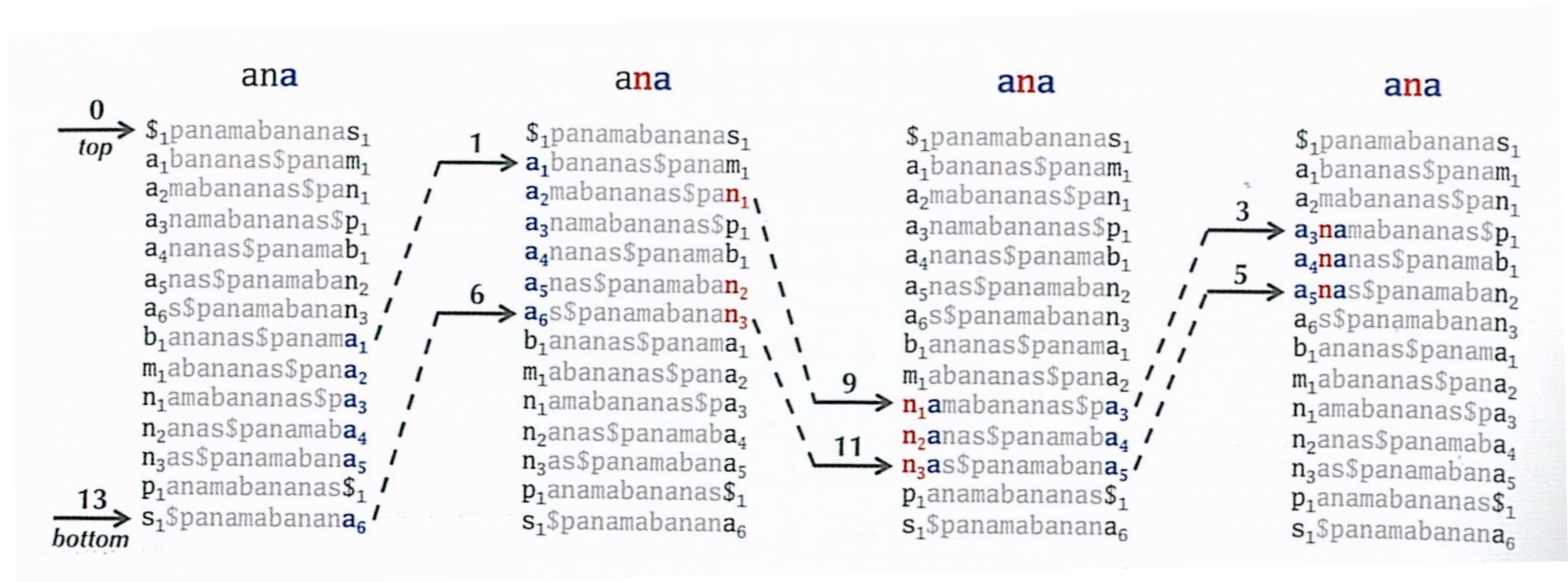# String to search **ana**

**Top-left matrix:**

```
$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
```

Application of **Rule 2**
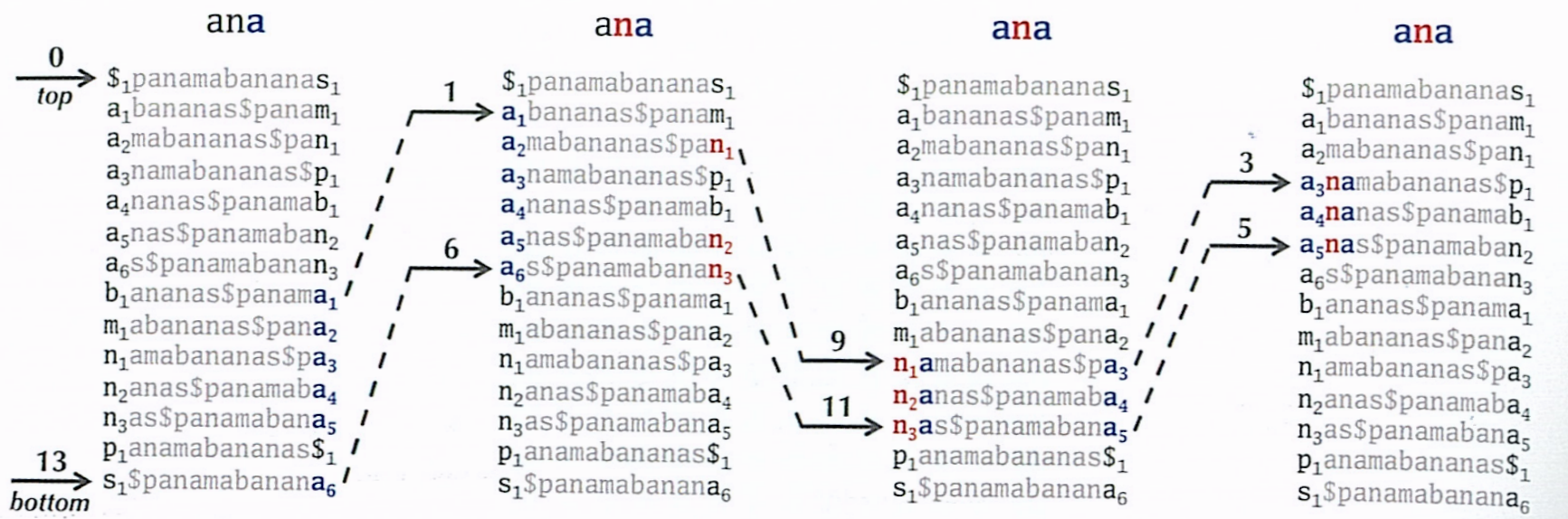
**Top-right matrix:**

```
$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 **a** are associated to an **n** in the **LastColumn**

All three **n** are associated to an **a** in the **FirstColumn**

**Bottom-left matrix:**

```
$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
```

**Bottom-right matrix:**

```
$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
```

# Pattern matching with BWT - Using pointers



The pointers top and bottom hold the indices of the first and last rows of M(*Text*) matching the current suffix of *Pattern*. The above diagram shows how these pointers are updated when walking backwards through **ana** and looking for substring matches in **panamabananas$**

**ana**      **ana**     **ana**     **ana**

Column 1 (top = 0, bottom = 13):

$\$_1$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$_2$
$a_6$\$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$_6$

Given a symbol at position i of **LastColumn**, the *Last-to-First* mapping identifies this symbol's position in **FirstColumn**

| i | FirstColumn | LastColumn | LASTTOFIRST(i) |
|---|---|---|---|
| 0 | $\$_1$ | $s_1$ | 13 |
| 1 | $a_1$ | $m_1$ | 8 |
| 2 | $a_2$ | $n_1$ | 9 |
| 3 | $a_3$ | $p_1$ | 12 |
| 4 | $a_4$ | $b_1$ | 7 |
| 5 | $a_5$ | $n_2$ | 10 |
| 6 | $a_6$ | $n_3$ | 11 |
| 7 | $b_1$ | $a_1$ | 1 |
| 8 | $m_1$ | $a_2$ | 2 |
| 9 | $n_1$ | $a_3$ | 3 |
| 10 | $n_2$ | $a_4$ | 4 |
| 11 | $n_3$ | $a_5$ | 5 |
| 12 | $p_1$ | $\$_1$ | 0 |
| 13 | $s_1$ | $a_6$ | 6 |

# Pattern matching with BWT

```
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
```

# Returning to Our Original Problem

- We need to look at INEXACT matching in order to find variants.

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

# Returning to Our Original Problem

- We need to look at INEXACT matching in order to find variants.

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

# Method 1: Seeding

- Say that *Pattern* appears in *Genome* with 1 mismatch:

|         |                          |
|---------|--------------------------|
| *Pattern* | act**t**ggct             |
| *Genome*  | …ggcac act**a**ggct cc… |

# Method 1: Seeding

- Say that *Pattern* appears in *Genome* with 1 mismatch:

  Pattern        a c t **t** g g c t

  Genome      … g g c a c a c t **a** g g c t c c …

- One of the substrings must match!

# Method 1: Seeding

- **Theorem:** If *Pattern* occurs in *Genome* with *d* mismatches, then we can divide *Pattern* into *d* + 1 "equal" pieces and find at least one exact match.

X**X**XXXXXX**X**XXX**X**XXXXXXX**X**XXXXX**X**XXXXXXXXX**X**XXX

X**X**XXXXXX**X**XXX**X**XXXXXXXX**X**XXXXXX**X**XXXXXXX**X**XXX

# Method 1: Seeding

- **Theorem:** If *Pattern* occurs in *Genome* with *d* mismatches, then we can divide *Pattern* into *d* + 1 "equal" pieces and find at least one exact match.

# Method 1: Seeding

- **Theorem:** If *Pattern* occurs in *Genome* with *d* mismatches, then we can divide *Pattern* into *d* + 1 "equal" pieces and find at least one exact match.

# Method 1: Seeding

- Say we are looking for at most $d$ mismatches.

- Divide each of our strings into $d + 1$ smaller pieces, called **seeds**.

# Method 2: BWT Saves the Day Again

- Recall: searching for an**a** 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**$_6$

# Method 2: BWT Saves the Day Again

- Recall: searching for a**na** in panamabananas

$\$_1$ p a n a m a b a n a n a $\mathbf{s}_1$
$\mathbf{a}_1$ b a n a n a s \$ p a n a $\mathbf{m}_1$
$\mathbf{a}_2$ m a b a n a n a s \$ p a $\mathbf{n}_1$
$\mathbf{a}_3$ n a m a b a n a n a s \$ $\mathbf{p}_1$
$\mathbf{a}_4$ n a n a s \$ p a n a m a $\mathbf{b}_1$
$\mathbf{a}_5$ n a s \$ p a n a m a b a $\mathbf{n}_2$
$\mathbf{a}_6$ s \$ p a n a m a b a n a $\mathbf{n}_3$
$\mathbf{b}_1$ a n a n a s \$ p a n a m $\mathbf{a}_1$
$\mathbf{m}_1$ a b a n a n a s \$ p a n $\mathbf{a}_2$
$\mathbf{n}_1$ a m a b a n a n a s \$ p $\mathbf{a}_3$
$\mathbf{n}_2$ a n a s \$ p a n a m a b $\mathbf{a}_4$
$\mathbf{n}_3$ a s \$ p a n a m a b a n $\mathbf{a}_5$
$\mathbf{p}_1$ a n a m a b a n a n a s $\$_1$
$\mathbf{s}_1$ \$ p a n a m a b a n a n $\mathbf{a}_6$

# Pattern Matching Using BWT

**ana**

| $i$ | FirstColumn | LastColumn | LASTTOFIRST($i$) |
|---|---|---|---|
| 0 | $\$_1$ | $s_1$ | 13 |
| 1 | $a_1$ | $m_1$ | 8 |
| 2 | $a_2$ | $n_1$ | 9 |
| 3 | $a_3$ | $p_1$ | 12 |
| 4 | $a_4$ | $b_1$ | 7 |
| 5 | $a_5$ | $n_2$ | 10 |
| 6 | $a_6$ | $n_3$ | 11 |
| 7 | $b_1$ | $a_1$ | 1 |
| 8 | $m_1$ | $a_2$ | 2 |
| 9 | $n_1$ | $a_3$ | 3 |
| 10 | $n_2$ | $a_4$ | 4 |
| 11 | $n_3$ | $a_5$ | 5 |
| 12 | $p_1$ | $\$_1$ | 0 |
| 13 | $s_1$ | $a_6$ | 6 |

To complete the pattern matching algorithm we need to answer

Where are the matched patterns?

and

How do we deal the mis-match tolerance?

# Where are the matched patterns?

The suffixes of **panamabananas$** that begin the rows of M(**panamabananas$**) are highlighted, and the suffixes beginning with **ana** are shown in green.

The suffix array records the starting position of each suffix in *Text*.

|  | M(*Text*) | SUFFIXARRAY(*Text*) |
|---|---|---|
| | $ p a n a m a b a n a n a s | 13 |
| | a b a n a n a s $ p a n a m | 5 |
| | 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 | 1 |
| | a n a n a s $ p a n a m a b | 7 |
| | a n a s $ p a n a m a b a n | 9 |
| | a s $ p a n a m a b a n a n | 11 |
| | b a n a n a s $ p a n a m a | 6 |
| | m a b a n a n a s $ p a n a | 4 |
| | n a m a b a n a n a s $ p a | 2 |
| | n a n a s $ p a n a m a b a | 8 |
| | n a s $ p a n a m a b a n a | 10 |
| | p a n a m a b a n a n a s $ | 0 |
| | s $ p a n a m a b a n a n a | 12 |

panamabananas$

$_1$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$_6$

panamabananas$

$_1$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$_6$

panamabananas$

$_1$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$_6$

Partial Suffix Array
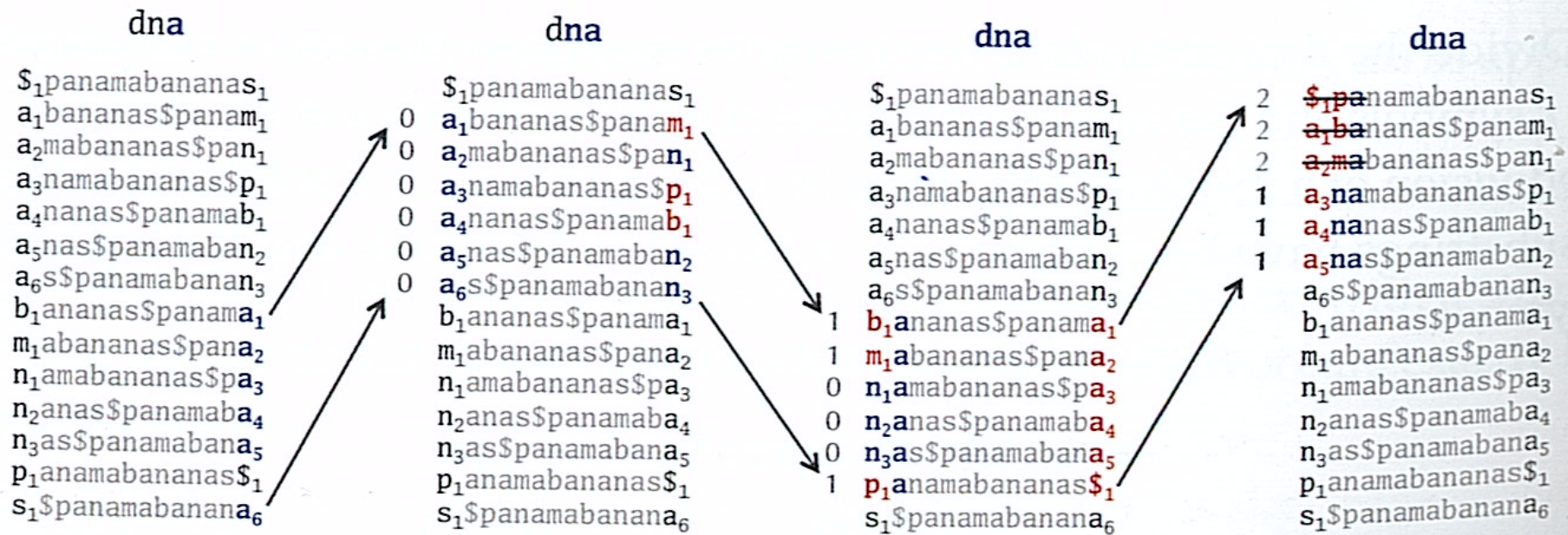
13
5
3
1
7
9
11
6
4
3
8
10
0
12

# How do we deal the mis-match tolerance?

Pattern **acttaggctcgggataatcc**
Text   **actaagtctcgggataagcc**

**Theorem**. If two strings of length n match with at most d mismatches, then they must share at least one k-mer of length k = [n/(d+1)]

To extend the BWT to approximate pattern matching, we will not stop when we encounter a mismatch. We proceed onward until we either find an approximate match or exceed the limit of $d$ mismatches.
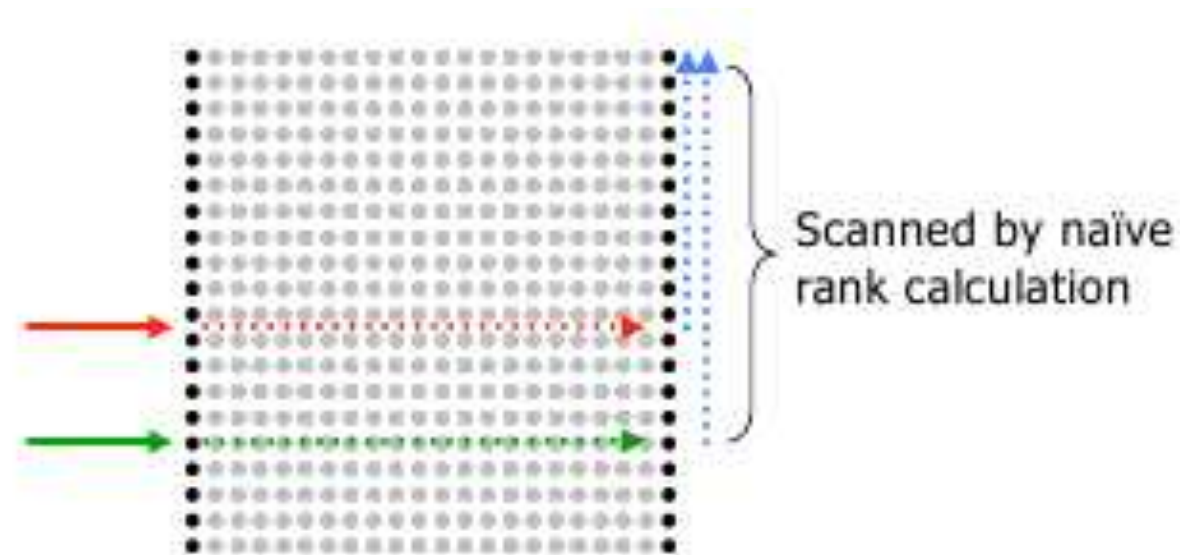
dna

```
$_1panamabananas_1
a_1bananas$panam_1
a_2mabananas$pan_1
a_3namabananas$p_1
a_4nanas$panamab_1
a_5nas$panamaban_2
a_6s$panamabanan_3
b_1ananas$panama_1
m_1abananas$pana_2
n_1amabananas$pa_3
n_2anas$panamaba_4
n_3as$panamabana_5
p_1anamabananas$_1
s_1$panamabanana_6
```

dna

```
  $_1panamabananas_1
0 a_1bananas$panam_1
0 a_2mabananas$pan_1
0 a_3namabananas$p_1
0 a_4nanas$panamab_1
0 a_5nas$panamaban_2
0 a_6s$panamabanan_3
  b_1ananas$panama_1
  m_1abananas$pana_2
  n_1amabananas$pa_3
  n_2anas$panamaba_4
  n_3as$panamabana_5
  p_1anamabananas$_1
  s_1$panamabanana_6
```

dna

```
  $_1panamabananas_1
  a_1bananas$panam_1
  a_2mabananas$pan_1
  a_3namabananas$p_1
  a_4nanas$panamab_1
  a_5nas$panamaban_2
  a_6s$panamabanan_3
1 b_1ananas$panama_1
1 m_1abananas$pana_2
0 n_1amabananas$pa_3
0 n_2anas$panamaba_4
0 n_3as$panamabana_5
1 p_1anamabananas$_1
  s_1$panamabanana_6
```

dna

```
2 $_1panamabananas_1
2 a_1bananas$panam_1
2 a_2mabananas$pan_1
1 a_3namabananas$p_1
1 a_4nanas$panamab_1
1 a_5nas$panamaban_2
  a_6s$panamabanan_3
  b_1ananas$panama_1
  m_1abananas$pana_2
  n_1amabananas$pa_3
  n_2anas$panamaba_4
  n_3as$panamabana_5
  p_1anamabananas$_1
  s_1$panamabanana_6
```

The number of mismatches encountered in a given row is shown in the column on the left.

# Key properties of Burrows-Wheeler Transform

- **Very little memory usage. Same as input (or less)**
  - Don't represent matrix, or strings, just pointers
  - Encode: Simply sort pointers. Decode: follow pointers
- **Original application: string compression (bZip2)**
  - Runs of letters compressed into (letter, runlength) pairs
- **Bioinformatics applications: substring searching**
  - Achieve similar run time as hash tables, suffix trees
  - But: very memory efficient ➔ practical speed gains
- **Mapping 100,000s of reads: only transform once**
  - Pre-process once; read counts in transformed space.
  - Reverse transform once, map counts to genome coords

- **LF**$(i, c)$ determines the rank of $qc$ in row $i$

- Nave way: count occurrences of $qc$ in all previous rows, complexity is linear in length of text  too slow



Scanned by naïve rank calculation

- Solution: pre-calculate cumulative counts for A/C/G/T up to periodic checkpoints in BWT

# Relationship Between BWT and Suffix Arrays

s = appellee$
123456789

| BWT matrix | The suffixes | Suffix array | subtract 1 | Suffix position - 1 |
|---|---|---|---|---|
| $appellee | $ | 9 | | s[9-1] = e |
| appellee$ | appellee$ | 1 | | s[1-1] = $ |
| e$appelle | e$ | 8 | | s[8-1] = e |
| ee$appell | ee$ | 7 | | s[7-1] = l |
| ellee$app | ellee$ | 4 | → | s[4-1] = p |
| lee$appel | lee$ | 6 | | s[6-1] = l |
| llee$appe | llee$ | 5 | | s[5-1] = e |
| pellee$ap | pellee$ | 3 | | s[3-1] = p |
| ppellee$a | ppellee$ | 2 | | s[2-1] = a |

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

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

| Suffixes | ID | Sorted Suffixes | Suffix Array | Sorted Rotations ($A_s$ matrix) | BWT Output ($L$) |
|---|---|---|---|---|---|
| mississippi$ | 1 | $ | 12 | $mississippi | i |
| ississippi$ | 2 | i$ | 11 | i$mississipp | p |
| ssissippi$ | 3 | ippi$ | 8 | ippi$mississ | s |
| sissippi$ | 4 | issippi$ | 5 | issippi$miss | s |
| issippi$ | 5 | ississippi$ | 2 | ississippi$m | m |
| ssippi$ | 6 | mississippi$ | 1 | mississippi$ | $ |
| sippi$ | 7 | pi$ | 10 | pi$mississip | p |
| ippi$ | 8 | ppi$ | 9 | ppi$mississi | i |
| ppi$ | 9 | sippi$ | 7 | sippi$missis | s |
| pi$ | 10 | sissippi$ | 4 | sissippi$mis | s |
| i$ | 11 | ssippi$ | 6 | ssippi$missi | i |
| $ | 12 | ssissippi$ | 3 | ssissippi$mi | i |

- Entire FM Index on DNA reference consists of:

  – BWT (same size as x)

  – Checkpoints ( 15% size of x)

  – SA sample ( 50% size of x)

- Total:  1.65x the size of x



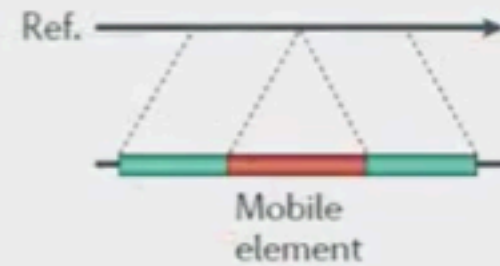~1.65x          >45x          >15x          >15x

# Structure Variations

### Deletion

Ref.

### Novel sequence insertion

Ref.

### Mobile-element insertion

Ref.

Mobile element

### Tandem duplication

Ref.

### Interspersed duplication

Ref.

### Inversion

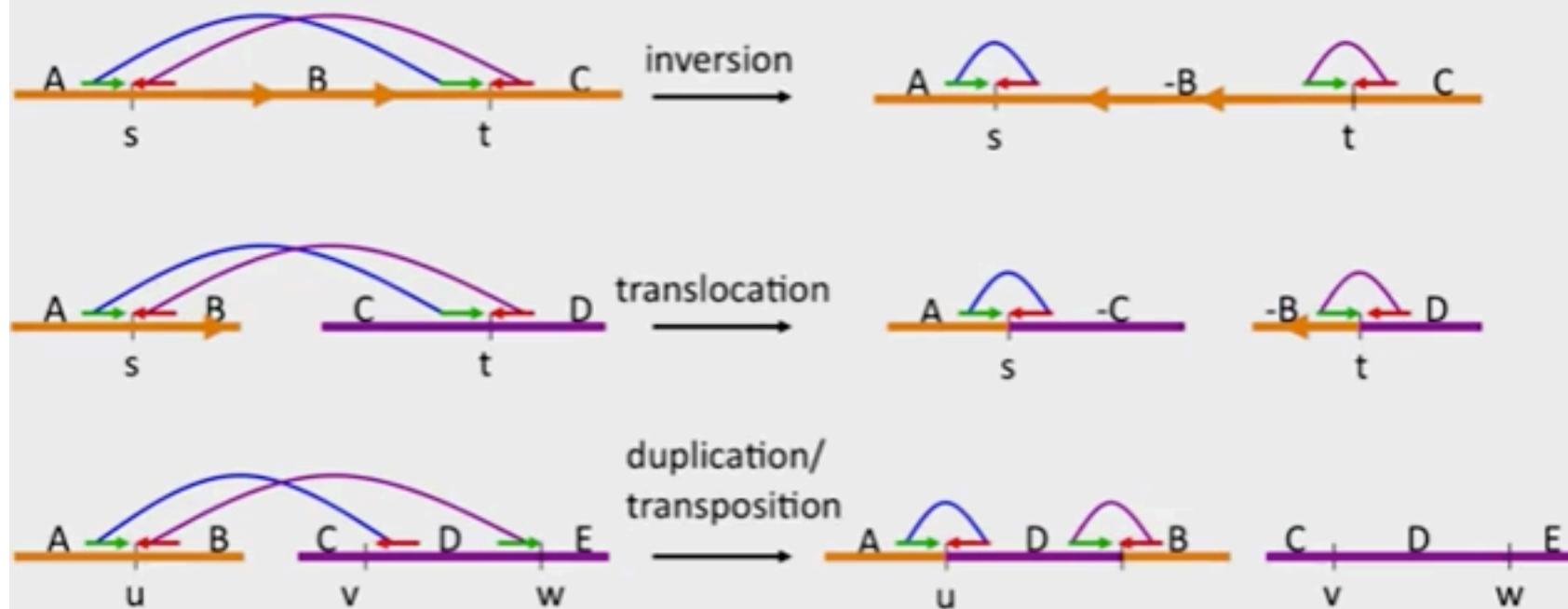Ref.

### Translocation

Ref.

Ref.
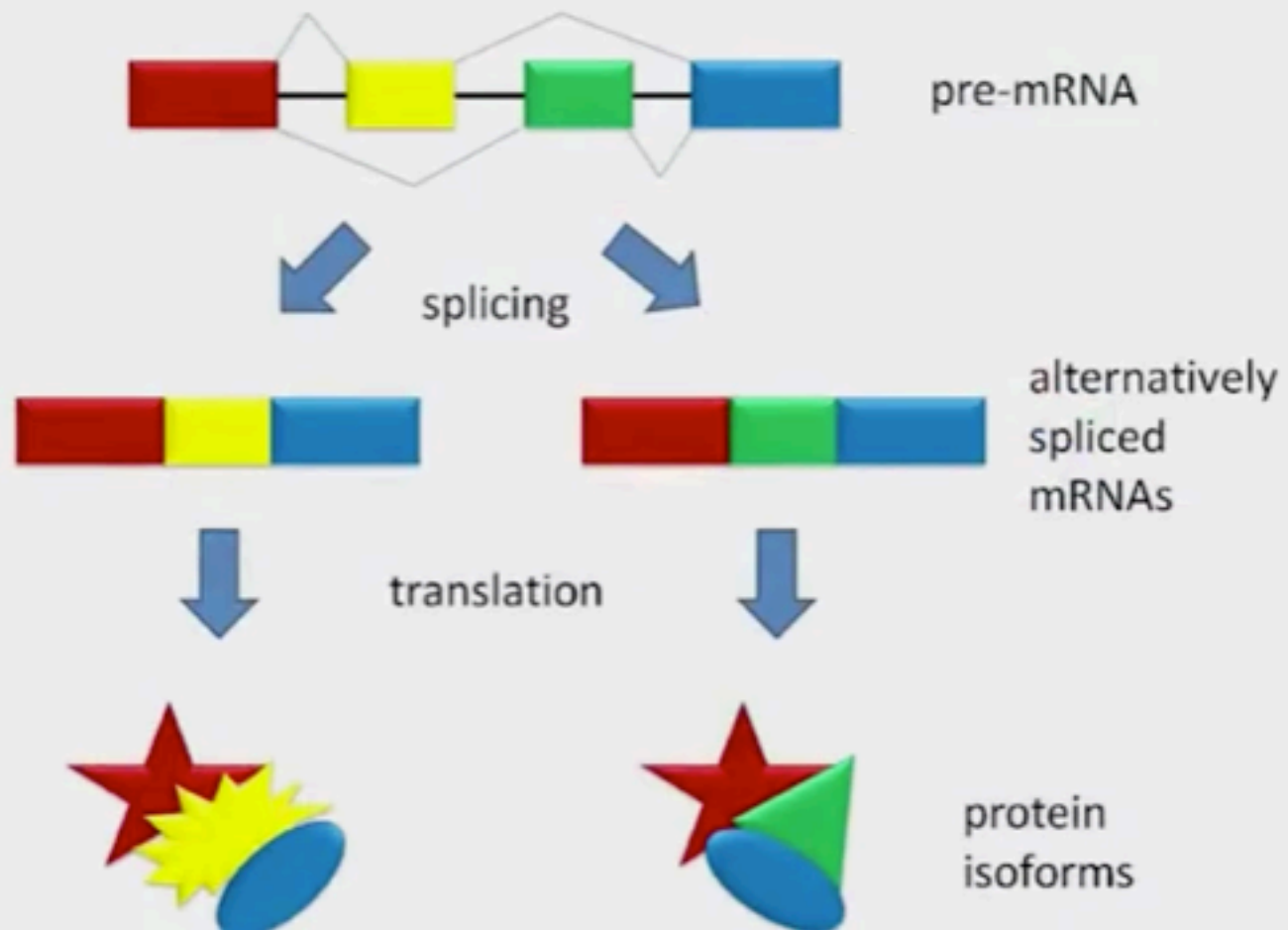
Alkan et al, 2012 Nature reviews

Variant Calling by Read Mapping
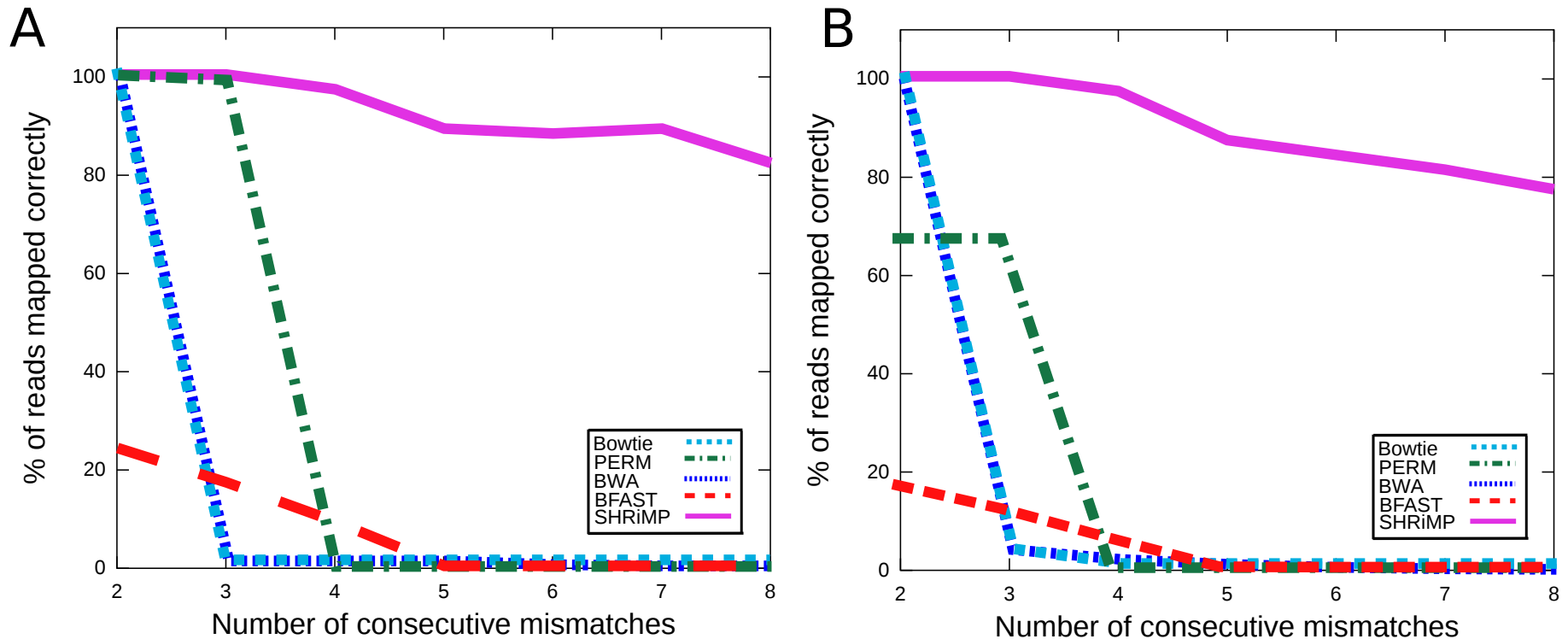
# Discovering Splice Junctions with RNA-Seq

# Algorithms comparison

| Tool | Space | Detection rate | | | | Time | Allowed mismatches-Seed |
|------|-------|---|---|------|------|--------|------------------------|
|      |       | 0 | 1 | 2 | 3 |        |                         |
| PerM | cs | 1 | 1 | 0.99 | 0.17 | 4m 30s | $F_4$ |
|      | nt | 1 | 1 | 1 | 0.81 | 3m |  |
| SOCS | cs | 1 | 1 | 0.18 | 0.02 | 50m | 3 |
| SHRiMP | cs | 1 | 1 | 1 | 0.69 | 125m | – |
|      | nt | 1 | 1 | 1 | 1 | 119m |  |
| MAQ | cs | 1 | 1 | 0.04 | 0 | 6m | 3 |
|     | nt | 1 | 1 | 1 | 0.01 | 6m |  |
| Bowtie | cs | 1 | 1 | 0.09 | 0.01 | 6m | 3 |
|        | nt | 1 | 1 | 0.84 | 0.0 | 5m 45s |  |

The accuracy of the various tools considering a synthetic dataset with continuous mismatches.

# Algorithms comparison



## Execution Time

| Program Name | 26 bp reads | 28 bp reads | 50 bp reads |
|:---|:---:|:---:|:---:|
| **SHRiMP** | 169.18 | 81.23 | 793.16 |
| **PERM** | 1.3 | 0.41 | 1 |
| **BWA** | 1.75 | 1.64 | 1.65 |
| **Bowtie** | 2.27 | 1.66 | 1.71 |
| **Bfast** | 11.44 | 11.32 | 11.98 |

Accuracy of the various tools over synthetic dataset and the execution time of each tool.