

# Linear-Time BWT Construction in Small Space

[Munro, Navarro & Nekrich, SODA 2017]

(slides by Patrick Dinklage, released under [CC0](#))

# Motivation

[Munro et al., SODA 2017]

# Motivation

[Munro et al., SODA 2017]

- $BWT[i] = T[SA[i]-1]$  well-known formula for  $O(n)$ -time *BWT*

# Motivation

[Munro et al., SODA 2017]

- $BWT[i] = T[SA[i]-1]$  well-known formula for  $O(n)$ -time *BWT*
- *BWT* can be stored in  $O(n \lg \sigma)$  bits (uncompressed)

# Motivation

[Munro et al., SODA 2017]

- $BWT[i] = T[SA[i]-1]$  well-known formula for  $O(n)$ -time *BWT*
- *BWT* can be stored in  $O(n \lg \sigma)$  bits (uncompressed)
- But: algorithms using *SA* require  $O(n \lg n)$  bits of working space!
  - ... or they no longer run in (deterministic) time  $O(n)$

# Motivation

[Munro et al., SODA 2017]

- $BWT[i] = T[SA[i]-1]$  well-known formula for  $O(n)$ -time  $BWT$
- $BWT$  can be stored in  $O(n \lg \sigma)$  bits (uncompressed)
- But: algorithms using  $SA$  require  $O(n \lg n)$  bits of working space!
  - ... or they no longer run in (deterministic) time  $O(n)$
- New  $O(n)$ -time algorithm using only  $O(n \lg \sigma)$  bits of space

# Motivation

[Munro et al., SODA 2017]

- $BWT[i] = T[SA[i]-1]$  well-known formula for  $O(n)$ -time  $BWT$
- $BWT$  can be stored in  $O(n \lg \sigma)$  bits (uncompressed)
- But: algorithms using  $SA$  require  $O(n \lg n)$  bits of working space!
  - ... or they no longer run in (deterministic) time  $O(n)$
- New  $O(n)$ -time algorithm using only  $O(n \lg \sigma)$  bits of space
- Sorting and batch operations – room for **efficient parallelization!**
  - [Fuentes-Sepúlveda et al., Theor. Comput. Sci. 812, 2020]

# Running Example

$$n = 30, \sigma = 5$$

123456789012345678901234567890  
*T* = GATCAATGAGGTGGACACCAGAGGCGGGG\$



# Running Example

$$n = 30, \sigma = 5$$

123456789012345678901234567890  
 $T = \text{GATCAATGAGGTGGACACCAGAGGCGGGG\$}$

Rotations:

$T_1 = \text{\$GATCAATGAGGTGGACACCAGAGGCGGGG}$

# Running Example

$$n = 30, \sigma = 5$$

123456789012345678901234567890  
 $T = \text{GATCAATGAGGTGGACACCAGAGGCGGGG\$}$

Rotations:

$T_1 = \text{\$GATCAATGAGGTGGACACCAGAGGCGGGG}$

$T_2 = \text{G\$GATCAATGAGGTGGACACCAGAGGCGGG}$

⋮

# Metasymbols

$$n = 30, \sigma = 5$$

$$\Delta := \lceil \lg_{\sigma}(n) \rceil = 3$$

# Metasymbols

$$n = 30, \sigma = 5$$

$$\Delta := \lceil \lg_{\sigma}(n) \rceil = 3$$

	1	2	3	4	5	6	7	8	9	0
$T_1 =$	\$GA	TCA	ATG	AGG	TGG	ACA	CCA	GAG	GCG	GGG
	1	2	3	4	5	6	7	8	9	0
$T_2 =$	G\$G	ATC	AAT	GAG	GTG	GAC	ACC	AGA	GGC	GGG

- › Combine  $\Delta$ -grams of  $T_1$  and  $T_2$  into metasymbols

# Metasymbols

$$n = 30, \sigma = 5$$

$$\Delta := \lceil \lg_{\sigma}(n) \rceil = 3$$

	1	2	3	4	5	6	7	8	9	0
$T_1 =$	\$GA	TCA	ATG	AGG	TGG	ACA	CCA	GAG	GCG	GGG
	1	2	3	4	5	6	7	8	9	0
$T_2 =$	G\$G	ATC	AAT	GAG	GTG	GAC	ACC	AGA	GGC	GGG

- Combine  $\Delta$ -grams of  $T_1$  and  $T_2$  into metasymbols
- Determine their lex. ranking in time  $O(n)$  (e.g., using radix sort)

# Metasymbols

$$n = 30, \sigma = 5$$

$$\Delta := \lceil \lg_{\sigma}(n) \rceil = 3$$

	1	2	3	4	5	6	7	8	9	0
$T_1 =$	\$GA	TCA	ATG	AGG	TGG	ACA	CCA	GAG	GCG	GGG
	1	2	3	4	5	6	7	8	9	0
$T_2 =$	G\$G	ATC	AAT	GAG	GTG	GAC	ACC	AGA	GGC	GGG

- Combine  $\Delta$ -grams of  $T_1$  and  $T_2$  into metasymbols
- Determine their lex. ranking in time  $O(n)$  (e.g., using radix sort)
- Compute  $SA(T_1T_2)$  in time  $O(n/\Delta)$ , space  $O((n/\Delta) \lg(n/\Delta)) = O(n \lg \sigma)$

# Metasymbols

$$n = 30, \sigma = 5$$

$$\Delta := \lceil \lg_{\sigma}(n) \rceil = 3$$

	1	2	3	4	5	6	7	8	9	0
$T_1 =$	\$GA	TCA	ATG	AGG	TGG	ACA	CCA	GAG	GCG	GGG
	1	2	3	4	5	6	7	8	9	0
$T_2 =$	G\$G	ATC	AAT	GAG	GTG	GAC	ACC	AGA	GGC	GGG

$$SA(T_1 T_2) = 1, 13, 6, 17, 18, 4, 12, 3, 7, 11, 16, 8, 14, 9, 19, 20, 10, 15, 2, 5$$

# Sequential Algorithm

- $SA(T_1T_2)$  equivalent to suffix array w.r.t. suffixes  $i\Delta$  and  $i\Delta-1$  of  $T$



# Sequential Algorithm

- $SA(T_1T_2)$  equivalent to suffix array w.r.t. suffixes  $i\Delta$  and  $i\Delta-1$  of  $T$

123456789012345678901234567890  
 $T = \text{GATCAATGAGGTGGACACCAGAGGC GGGS}$

# Sequential Algorithm

- $SA(T_1T_2)$  equivalent to suffix array w.r.t. suffixes  $i\Delta$  and  $i\Delta-1$  of  $T$

123456789012345678901234567890  
 $T = \text{GATCAATGAGGTGGACACCAGAGGCGGGG\$}$

- Append symbols at  $i\Delta-1$  and  $i\Delta-2$  to initially empty string  $B$ 
  - ... via one scan of  $SA(T_1T_2)$
  - In  $B$ , we will construct the BWT of  $T$  step by step

# Example: Sequential Algorithm

$SA(T_1T_2) = 1, 13, 6, 17, 18, 4, 12, 3, 7, 11, 16, 8, 14, 9, 19, 20, 10, 15, 2, 5$

	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>	<b>0</b>
$T_1 =$	<b>\$GA</b>	TCA	ATG	AGG	TGG	ACA	CCA	GAG	GCG	GGG

	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>	<b>0</b>
$T_2 =$	G <b>\$G</b>	ATC	AAT	GAG	GTG	GAC	ACC	AGA	GGC	GGG

123456789012345678901234567890  
 $T =$  GATCAATGAGGTGGACACCAGAGGGCGGG**G\$**

# Example: Sequential Algorithm

$SA(T_1T_2) = 1, 13, 6, 17, 18, 4, 12, 3, 7, 11, 16, 8, 14, 9, 19, 20, 10, 15, 2, 5$



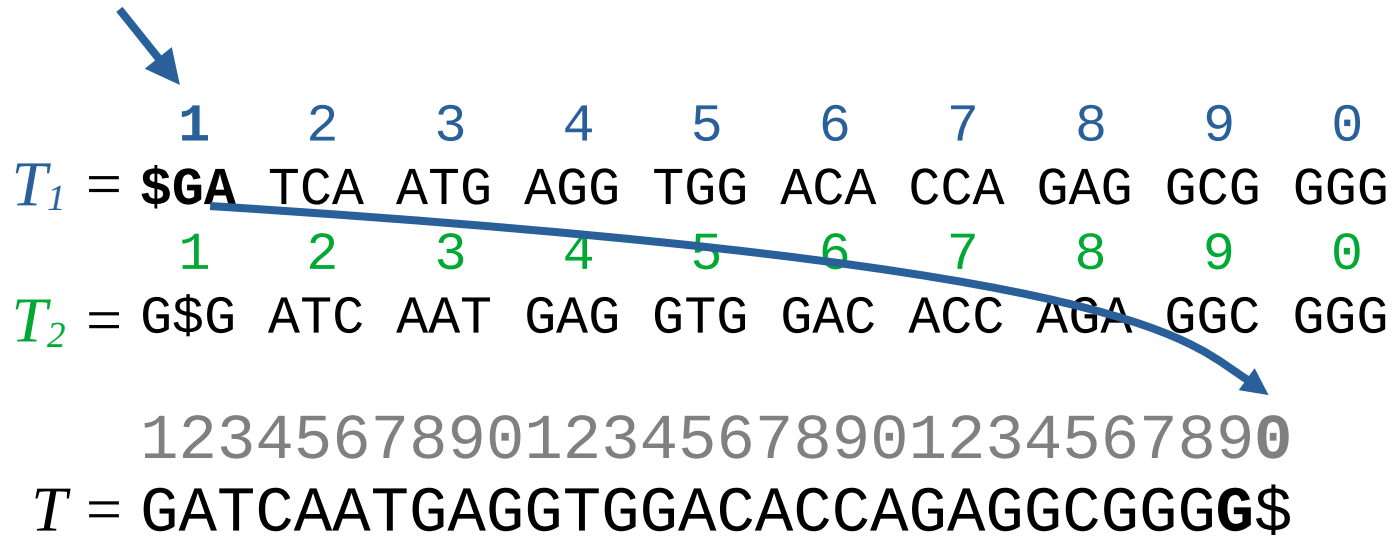
$T_1 =$  **\$GA** TCA ATG AGG TGG ACA CCA GAG GCG GGG

$T_2 =$  G**\$G** ATC AAT GAG GTG GAC ACC AGA GGC GGG

123456789012345678901234567890  
 $T =$  GATCAATGAGGTGGACACCAGAGGC**GGGG**\$

# Example: Sequential Algorithm

$SA(T_1T_2) = 1, 13, 6, 17, 18, 4, 12, 3, 7, 11, 16, 8, 14, 9, 19, 20, 10, 15, 2, 5$



# Example: Sequential Algorithm

$SA(T_1T_2) = 1, 13, 6, 17, 18, 4, 12, 3, 7, 11, 16, 8, 14, 9, 19, 20, 10, 15, 2, 5$

$T_1 =$   $\$GA$  TCA ATG AGG TGG ACA CCA GAG GCG GGG  
          1   2   3   4   5   6   7   8   9   0  
 $T_2 =$  G\$G ATC AAT GAG GTG GAC ACC AGA GGC GGG

123456789012345678901234567890  
 $T =$  GATCAATGAGGTGGACACCAGAGGCGGGG\$

$B = G$

# Example: Sequential Algorithm

$SA(T_1T_2) = 1, 13, 6, 17, 18, 4, 12, 3, 7, 11, 16, 8, 14, 9, 19, 20, 10, 15, 2, 5$

	1	2	3	4	5	6	7	8	9	0
$T_1$	= \$GA	TCA	ATG	AGG	TGG	ACA	CCA	GAG	GCG	GGG
	1	2	3	4	5	6	7	8	9	0
$T_2$	= G\$G	ATC	<b>AAT</b>	GAG	GTG	GAC	ACC	AGA	GGC	GGG

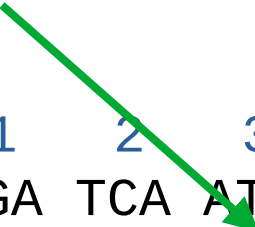
123456789012345678901234567890  
 $T = \text{GATCAATGAGGTGGACACCAGAGGCGGGG\$}$

$B = \text{GC}$

# Example: Sequential Algorithm

$SA(T_1T_2) = 1, 13, 6, 17, 18, 4, 12, 3, 7, 11, 16, 8, 14, 9, 19, 20, 10, 15, 2, 5$

	1	2	3	4	5	6	7	8	9	0	
$T_1$	=	\$GA	TCA	ATG	AGG	TGG	ACA	CCA	GAG	GCG	GGG
		1	2	3	4	5	6	7	8	9	0
$T_2$	=	G\$G	ATC	AAT	GAG	GTG	GAC	ACC	AGA	GGC	GGG



123456789012345678901234567890  
 $T = \text{GATCAATGAGGTGGACACCAGAGGCGGGG\$}$

$B = \text{GC}$



# Example: Sequential Algorithm

$SA(T_1T_2) = 1, 13, 6, 17, 18, 4, 12, 3, 7, 11, 16, 8, 14, 9, 19, 20, 10, 15, 2, 5$

$T_1 =$ 

	1	2	3	4	5	6	7	8	9	0
\$GA	TCA	ATG	AGG	TGG	ACA	CCA	GAG	GCG	GGG	

$T_2 =$ 

	1	2	3	4	5	6	7	8	9	0
G\$G	ATC	AAT	GAG	GTG	GAC	ACC	AGA	GGC	GGG	

	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	
$T =$	G	A	<b>C</b>	A	A	T	G	A	G	G	T	G	G	A	C	A	C	C	A	G	A	G	G	C	G	G	G	G	G	G	\$

$B = GC$

# Example: Sequential Algorithm

$SA(T_1T_2) = 1, 13, 6, 17, 18, 4, 12, 3, 7, 11, 16, 8, 14, 9, 19, 20, 10, 15, 2, 5$

	1	2	3	4	5	6	7	8	9	0	
$T_1$	=	\$GA	TCA	ATG	AGG	TGG	ACA	CCA	GAG	GCG	GGG
		1	2	3	4	5	6	7	8	9	0
$T_2$	=	G\$G	ATC	AAT	GAG	GTG	GAC	ACC	AGA	GGC	GGG

123456789012345678901234567890  
 $T = \text{GATCAATGAGGTGGACACCAGAGGC GGGG\$}$

$B = \text{GC}$

# Example: Sequential Algorithm

$SA(T_1T_2) = 1, 13, 6, 17, 18, 4, 12, 3, 7, 11, 16, 8, 14, 9, 19, 20, 10, 15, 2, 5$

	1	2	3	4	5	6	7	8	9	0
$T_1$	= \$GA	TCA	ATG	AGG	TGG	ACA	CCA	GAG	GCG	GGG
	1	2	3	4	5	6	7	8	9	0
$T_2$	= G\$G	ATC	AAT	GAG	GTG	GAC	ACC	AGA	GGC	GGG

123456789012345678901234567890  
 $T = \text{GATCAATGAGGTGGACACCAGAGGCGGGG\$}$

$B = \text{GCGCCGGAAGGATGACGGAG}$

# Grand Strategy

# Grand Strategy

- Build BWT of  $T$  in a total of  $\Delta$  steps
  - Each step will take time  $O(n/\Delta)$ , so we need time  $O(n)$  in total

# Grand Strategy

- Build BWT of  $T$  in a total of  $\Delta$  steps
  - Each step will take time  $O(n/\Delta)$ , so we need time  $O(n)$  in total
- After step  $j$ ,  $B$  contains the BWT of  $T$  with respect to suffixes at positions  $i\Delta-j, \dots, i\Delta-1$ 
  - Let  $T[k\dots]$  be the  $(i+1)$ -th smallest suffix, then  $B[k] := T[(k-1) \bmod n]$
  - $B$  is a substring of the BWT of  $T$

# Grand Strategy

- Build BWT of  $T$  in a total of  $\Delta$  steps
  - Each step will take time  $O(n/\Delta)$ , so we need time  $O(n)$  in total
- After step  $j$ ,  $\mathbf{B}$  contains the BWT of  $T$  with respect to suffixes at positions  $i\Delta-j, \dots, i\Delta-1$ 
  - Let  $T[k\dots]$  be the  $(i+1)$ -th smallest suffix, then  $B[k] := T[(k-1) \bmod n]$
  - $\mathbf{B}$  is a substring of the BWT of  $T$
- In each step, we insert  $n/\Delta$  further characters into  $\mathbf{B}$

# Grand Strategy

- Build BWT of  $T$  in a total of  $\Delta$  steps
  - Each step will take time  $O(n/\Delta)$ , so we need time  $O(n)$  in total
- After step  $j$ ,  $\mathbf{B}$  contains the BWT of  $T$  with respect to suffixes at positions  $i\Delta-j, \dots, i\Delta-1$ 
  - Let  $T[k\dots]$  be the  $(i+1)$ -th smallest suffix, then  $B[k] := T[(k-1) \bmod n]$
  - $\mathbf{B}$  is a substring of the BWT of  $T$
- In each step, we insert  $n/\Delta$  further characters into  $\mathbf{B}$
- Steps 1 and 2 were already combined



# Preparation for Step $j+1$

## Preparation for Step $j+1$

- While computing  $B$ , also build two arrays for further steps:

## Preparation for Step $j+1$

- While computing  $B$ , also build two arrays for further steps:
- For every suffix of  $T_j$  (metacharacters), store the  $B$ -position of its first character in array  $W$
- Space  $O((n/\Delta) \lg n) = O(n \lg \sigma)$

# Preparation for Step $j+1$

- While computing  $B$ , also build two arrays for further steps:
- For every suffix of  $T_j$  (metacharacters), store the  $B$ -position of its first character in array  $W$ 
  - Space  $O((n/\Delta) \lg n) = O(n \lg \sigma)$
- $C$  array of  $B$  (for every character, # of lex. smaller characters)
  - Space  $O(\sigma \lg n) = O(n \lg \sigma)$  if  $\sigma = o(n)$

## Preparation for Step 3

- For every suffix of  $T_2$  (metacharacters), store the  $B$ -position of its first character in array  $W$

## Preparation for Step 3

- For every suffix of  $T_2$  (metacharacters), store the  $B$ -position of its first character in array  $W$

$SA(T_1T_2) =$  1, 13, 6, 17, 18, 4, 12, 3, 7, 11, 16, 8, 14, 9, 19, 20, 10, 15, 2, 5

$B =$  GCGCCGGAAGGATGACGGAG

## Preparation for Step 3

- For every suffix of  $T_2$  (metacharacters), store the  $B$ -position of its first character in array  $W$

1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0  
 $SA(T_1T_2) = 1, 13, 6, 17, 18, 4, 12, 3, 7, 11, 16, 8, 14, 9, 19, 20, 10, 15, 2, 5$

$B = \text{GCGCCGGAAGGATGACGGAG}$

$W[1]=10, W[2]=7, W[3]=2, W[4]=13, W[5]=18$   
 $W[6]=11, W[7]=4, W[8]=5, W[9]=15, W[10]=16$

## Preparation for Step 3

- For every suffix of  $T_2$  (metacharacters), store the  $B$ -position of its first character in array  $W$

SA( $T_1T_2$ ) = 1, 13, 6, 17, 18, 4, 12, 3, 7, 11, 16, 8, 14, 9, 19, 20, 10, 15, 2, 5

$B =$  GCGCCGGAAGGATGACGGAG

$W[1]=10, W[2]=7, W[3]=2, W[4]=13, W[5]=18$   
 $W[6]=11, W[7]=4, W[8]=5, W[9]=15, W[10]=16$



## Preparation for Step 3

- For every suffix of  $T_2$  (metacharacters), store the  $B$ -position of its first character in array  $W$

SA( $T_1T_2$ ) = 1, 13, 6, 17, 18, 4, 12, 3, 7, 11, 16, 8, 14, 9, 19, 20, 10, 15, 2, 5



$B = \text{GCGCCGGAAGGATGACGGAG}$

$W[1]=10, W[2]=7, W[3]=2, W[4]=13, W[5]=18$   
 $W[6]=11, W[7]=4, W[8]=5, W[9]=15, W[10]=16$

## Preparation for Step 3

- For every suffix of  $T_2$  (metacharacters), store the  $B$ -position of its first character in array  $W$

$B = \text{GCGCCGGAAGGATGACGGAG}$

$W[1]=10, W[2]=7, W[3]=2, W[4]=13, W[5]=18$   
 $W[6]=11, W[7]=4, W[8]=5, W[9]=15, W[10]=16$

$C[\$]=0, C[A]=0, C[C]=5, C[G]=9, C[T]=19$

Step  $j > 2$

## Step $j > 2$

- After step  $j-1$ ,  $B$  contains the BWT of  $T$  w.r.t. suffixes at positions  $i\Delta-j+1, \dots, i\Delta-1$

## Step $j > 2$

- After step  $j-1$ ,  $B$  contains the BWT of  $T$  w.r.t. suffixes at positions  $i\Delta-j+1, \dots, i\Delta-1$
- In step  $j$ , we insert suffixes  $i\Delta-j$  into  $B$ 
  - But where?

Step  $j > 2$

## Step $j > 2$

- › Let  $S_{j-1}$  be the set of suffixes of  $T$  starting at positions  $i\Delta-j+1$  to  $i\Delta-2$ 
  - › These have been inserted into  $B$  in steps 2 to  $j-1$
  - ›  $W$  contains the positions of suffixes  $i\Delta-j+1$  in  $B$

## Step $j > 2$

- Let  $S_{j-1}$  be the set of suffixes of  $T$  starting at positions  $i\Delta-j+1$  to  $i\Delta-2$ 
  - These have been inserted into  $B$  in steps 2 to  $j-1$
  - $W$  contains the positions of suffixes  $i\Delta-j+1$  in  $B$
- Let  $s$  be a suffix of  $T$  at position  $i\Delta-j$  starting with character  $\alpha$



## Step $j > 2$

- Let  $S_{j-1}$  be the set of suffixes of  $T$  starting at positions  $i\Delta-j+1$  to  $i\Delta-2$ 
  - These have been inserted into  $B$  in steps 2 to  $j-1$
  - $W$  contains the positions of suffixes  $i\Delta-j+1$  in  $B$
- Let  $s$  be a suffix of  $T$  at position  $i\Delta-j$  starting with character  $\alpha$
- For every occurrence of some  $\alpha' < \alpha$ , there is exactly one suffix  $s' < s$  starting with  $\alpha'$  in  $S_{j-1}$  (that's  $C[\alpha]$  in total)
  - Because we begin with the  $\Delta$ -th rotation of  $T$  and rotate to the right in each step

## Step $j > 2$

- Let  $S_{j-1}$  be the set of suffixes of  $T$  starting at positions  $i\Delta-j+1$  to  $i\Delta-2$ 
  - These have been inserted into  $B$  in steps 2 to  $j-1$
  - $W$  contains the positions of suffixes  $i\Delta-j+1$  in  $B$
- Let  $s$  be a suffix of  $T$  at position  $i\Delta-j$  starting with character  $\alpha$
- For every occurrence of some  $\alpha' < \alpha$ , there is exactly one suffix  $s' < s$  starting with  $\alpha'$  in  $S_{j-1}$  (that's  $C[\alpha]$  in total)
  - Because we begin with the  $\Delta$ -th rotation of  $T$  and rotate to the right in each step
- There are exactly  $\text{rank}_\alpha(B, W[i])$  suffixes  $s'' \leq s$  in  $S_{j-1}$  that start with  $\alpha$ 
  - Because the BWT is order-preserving

## Step $j > 2$

- Let  $S_{j-1}$  be the set of suffixes of  $T$  starting at positions  $i\Delta-j+1$  to  $i\Delta-2$ 
  - These have been inserted into  $B$  in steps 2 to  $j-1$
  - $W$  contains the positions of suffixes  $i\Delta-j+1$  in  $B$
- Let  $s$  be a suffix of  $T$  at position  $i\Delta-j$  starting with character  $\alpha$
- For every occurrence of some  $\alpha' < \alpha$ , there is exactly one suffix  $s' < s$  starting with  $\alpha'$  in  $S_{j-1}$  (that's  $C[\alpha]$  in total)
  - Because we begin with the  $\Delta$ -th rotation of  $T$  and rotate to the right in each step
- There are exactly  $\text{rank}_\alpha(B, W[i])$  suffixes  $s'' \leq s$  in  $S_{j-1}$  that start with  $\alpha$ 
  - Because the BWT is order-preserving
- Thus,  $C[\alpha] + \text{rank}_\alpha(B, W[i])$  suffixes from  $S_{j-1}$  precede  $s$  in  $B$

Step  $j > 2$

## Step $j > 2$

- › ... Thus,  $C[\alpha] + \text{rank}_\alpha(\mathbf{B}, \mathbf{W}[i])$  suffixes from  $S_{j-1}$  precede  $s$  in  $\mathbf{B}$

## Step $j > 2$

- › ... Thus,  $C[\alpha] + \text{rank}_\alpha(\mathbf{B}, \mathbf{W}[i])$  suffixes from  $S_{j-1}$  precede  $s$  in  $\mathbf{B}$
- › Consider  $SA(T_1 T_j)$  and let  $c_i$  be the # of  $T_1$  suffixes appearing before suffix  $i\Delta - j$  in  $SA(T_1 T_j)$

## Step $j > 2$

- ... Thus,  $C[\alpha] + \text{rank}_\alpha(\mathbf{B}, \mathbf{W}[i])$  suffixes from  $S_{j-1}$  precede  $s$  in  $\mathbf{B}$
- Consider  $SA(T_1 T_j)$  and let  $c_i$  be the # of  $T_1$  suffixes appearing before suffix  $i\Delta - j$  in  $SA(T_1 T_j)$
- These also precede  $s$  in  $\mathbf{B}$  and have not yet been considered

## Step $j > 2$

- ... Thus,  $C[\alpha] + \text{rank}_\alpha(\mathbf{B}, \mathbf{W}[i])$  suffixes from  $S_{j-1}$  precede  $s$  in  $\mathbf{B}$
- Consider  $SA(T_1 T_j)$  and let  $c_i$  be the # of  $T_1$  suffixes appearing before suffix  $i\Delta - j$  in  $SA(T_1 T_j)$
- These also precede  $s$  in  $\mathbf{B}$  and have not yet been considered
- In step  $j$ , we insert  $T[i\Delta - j]$  at position  $C[\alpha] + \text{rank}_\alpha(\mathbf{B}, \mathbf{W}[i]) + c_i$



### Example: Step 3

123456789012345678901234567890  
 $T = \text{GATCAATGAGGTGGACACCAGAGGGCGGGG}\$$

1 2 3 4 5 6 7 8 9 0  
 $T_1 = \text{\$GA TCA ATG AGG TGG ACA CCA GAG GCG GGG}$

1 2 3 4 5 6 7 8 9 0  
 $T_3 = \text{GG\$ GAT CAA TGA GGT GGA CAC CAG AGG CGG}$

$SA(T_1 T_3) = 1, 6, 19, 4, 3, 13, 17, 18, 7, 20, 8, 12, 9, 11, 16, 10, 15, 2, 14, 5$

12345678901234567890  
 $B = \text{GCGCCGGAAGGATGACGGAG}$

### Example: Step 3

123456789012345678901234567890  
 $T = \text{GATCAATGAGGTGGACACCAGAGGCGGG\$}$

1 2 3 4 5 6 7 8 9 0  
 $T_1 = \text{\$GA TCA ATG AGG TGG ACA CCA GAG GCG GGG}$

1 2 3 4 5 6 7 8 9 0  
 $T_3 = \text{GG\$ GAT CAA TGA GGT GGA CAC CAG AGG CGG}$

$SA(T_1 T_3) = 1, 6, 19, 4, 3, 13, 17, 18, 7, 20, 8, 12, 9, 11, 16, 10, 15, 2, 14, 5$

12345678901234567890  
 $B = \text{GCGCCGGAAGGATGACGGAG}$

- Insert **G** at position  $C[\text{G}] + \text{rank}_{\text{G}}(B, W[1]) + c_1 = 9 + 5 + 7 = 21$

## Example: Step 3

123456789012345678901234567890  
 $T = \mathbf{G}$ ATCAATGAGGTGGACACCAGAGGGCGGGG\$

1 2 3 4 5 6 7 8 9 0  
 $T_1 =$  \$GA TCA ATG AGG TGG ACA CCA GAG GCG GGG

1 2 3 4 5 6 7 8 9 0  
 $T_3 =$  GG\$  $\mathbf{G}$ AT CAA TGA GGT GGA CAC CAG AGG CGG

$SA(T_1 T_3) = 1, 6, 19, 4, 3, 13, 17, 18, 7, 20, 8, 12, 9, 11, 16, 10, 15, 2, 14, 5$

12345678901234567890  
 $B =$  GCGCCGGAAGGATGACGGAG

- Insert \$ at position  $C[\mathbf{G}] + \text{rank}_{\mathbf{G}}(B, W[2])^{=1} + c_2 = 9 + 4 + 6 = 19$

## Example: Step 3

12345678901234567890

**B** = GCGCCGGAAGGATGACGGAG

- Insert **G** at position  $C[G] + \text{rank}_G(B, W[1]) + c_1 = 9 + 5 + 7 = 21$
- Insert **\$** at position  $C[G] + \text{rank}_G(B, W[2]) + c_2 = 9 + 4 + 6 = 19$
- Insert **T** at position  $C[C] + \text{rank}_C(B, W[3]) + c_3 = 5 + 1 + 4 = 10$
- Insert **A** at position  $C[T] + \text{rank}_T(B, W[4]) + c_4 = 19 + 1 + 9 = 29$
- Insert **A** at position  $C[G] + \text{rank}_G(B, W[5]) + c_5 = 9 + 9 + 8 = 26$
- Insert **T** at position  $C[G] + \text{rank}_G(B, W[6]) + c_6 = 9 + 6 + 7 = 22$
- Insert **A** at position  $C[C] + \text{rank}_C(B, W[7]) + c_7 = 5 + 2 + 4 = 11$
- Insert **C** at position  $C[C] + \text{rank}_C(B, W[8]) + c_8 = 5 + 3 + 4 = 12$
- Insert **G** at position  $C[A] + \text{rank}_A(B, W[9]) + c_9 = 0 + 4 + 2 = 6$
- Insert **G** at position  $C[C] + \text{rank}_C(B, W[10]) + c_{10} = 5 + 4 + 5 = 14$

## Example: Step 3

123456789012345678901234567890

**B** = GCGCC - GGA - - - A - GGAT - G - - AGC - GA - G

- Insert **G** at position  $C[G] + \text{rank}_G(\mathbf{B}, \mathbf{W}[1]) + c_1 = 9 + 5 + 7 = 21$
- Insert **\$** at position  $C[G] + \text{rank}_G(\mathbf{B}, \mathbf{W}[2]) + c_2 = 9 + 4 + 6 = 19$
- Insert **T** at position  $C[C] + \text{rank}_C(\mathbf{B}, \mathbf{W}[3]) + c_3 = 5 + 1 + 4 = 10$
- Insert **A** at position  $C[T] + \text{rank}_T(\mathbf{B}, \mathbf{W}[4]) + c_4 = 19 + 1 + 9 = 29$
- Insert **A** at position  $C[G] + \text{rank}_G(\mathbf{B}, \mathbf{W}[5]) + c_5 = 9 + 9 + 8 = 26$
- Insert **T** at position  $C[G] + \text{rank}_G(\mathbf{B}, \mathbf{W}[6]) + c_6 = 9 + 6 + 7 = 22$
- Insert **A** at position  $C[C] + \text{rank}_C(\mathbf{B}, \mathbf{W}[7]) + c_7 = 5 + 2 + 4 = 11$
- Insert **C** at position  $C[C] + \text{rank}_C(\mathbf{B}, \mathbf{W}[8]) + c_8 = 5 + 3 + 4 = 12$
- Insert **G** at position  $C[A] + \text{rank}_A(\mathbf{B}, \mathbf{W}[9]) + c_9 = 0 + 4 + 2 = 6$
- Insert **G** at position  $C[C] + \text{rank}_C(\mathbf{B}, \mathbf{W}[10]) + c_{10} = 5 + 4 + 5 = 14$

## Example: Step 3

123456789012345678901234567890

**B** = GCGCCGGGATACAGGGAT\$GGTAGCAGAAG

- Insert **G** at position  $C[G] + \text{rank}_G(\mathbf{B}, \mathbf{W}[1]) + c_1 = 9 + 5 + 7 = 21$
- Insert **\$** at position  $C[G] + \text{rank}_G(\mathbf{B}, \mathbf{W}[2]) + c_2 = 9 + 4 + 6 = 19$
- Insert **T** at position  $C[C] + \text{rank}_C(\mathbf{B}, \mathbf{W}[3]) + c_3 = 5 + 1 + 4 = 10$
- Insert **A** at position  $C[T] + \text{rank}_T(\mathbf{B}, \mathbf{W}[4]) + c_4 = 19 + 1 + 9 = 29$
- Insert **A** at position  $C[G] + \text{rank}_G(\mathbf{B}, \mathbf{W}[5]) + c_5 = 9 + 9 + 8 = 26$
- Insert **T** at position  $C[G] + \text{rank}_G(\mathbf{B}, \mathbf{W}[6]) + c_6 = 9 + 6 + 7 = 22$
- Insert **A** at position  $C[C] + \text{rank}_C(\mathbf{B}, \mathbf{W}[7]) + c_7 = 5 + 2 + 4 = 11$
- Insert **C** at position  $C[C] + \text{rank}_C(\mathbf{B}, \mathbf{W}[8]) + c_8 = 5 + 3 + 4 = 12$
- Insert **G** at position  $C[A] + \text{rank}_A(\mathbf{B}, \mathbf{W}[9]) + c_9 = 0 + 4 + 2 = 6$
- Insert **G** at position  $C[C] + \text{rank}_C(\mathbf{B}, \mathbf{W}[10]) + c_{10} = 5 + 4 + 5 = 14$

## Example: Step 3

123456789012345678901234567890

**B** = GCGCCGGGATACAGGGAT\$GGTAGCAGAAG

- Insert **G** at position  $C[G] + \text{rank}_G(B, W[1]) + c_1 = 9 + 5 + 7 = 21$
- Insert **\$** at position  $C[G] + \text{rank}_G(B, W[2]) + c_2 = 9 + 4 + 6 = 19$
- Insert **T** at position  $C[C] + \text{rank}_C(B, W[3]) + c_3 = 5 + 1 + 4 = 10$
- Insert **A** at position  $C[T] + \text{rank}_T(B, W[4]) + c_4 = 19 + 1 + 9 = 29$
- Insert **A** at position  $C[G] + \text{rank}_G(B, W[5]) + c_5 = 9 + 9 + 8 = 26$
- Insert **T** at position  $C[G] + \text{rank}_G(B, W[6]) + c_6 = 9 + 6 + 7 = 22$
- Insert **A** at position  $C[C] + \text{rank}_C(B, W[7]) + c_7 = 5 + 2 + 4 = 11$
- Insert **C** at position  $C[C] + \text{rank}_C(B, W[8]) + c_8 = 5 + 3 + 4 = 12$
- Insert **G** at position  $C[A] + \text{rank}_A(B, W[9]) + c_9 = 0 + 4 + 2 = 6$
- Insert **G** at position  $C[C] + \text{rank}_C(B, W[10]) + c_{10} = 5 + 4 + 5 = 14$

$n/\Delta$  rank queries **must not exceed  $O(1)$  time per query (amortized)!**

# Data Structure for Batched Rank Queries

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

***B*** = G C G C C G G A A G G A T G A C G G A G



# Data Structure for Batched Rank Queries

$B = \begin{array}{|cccc|cccc|cccc|cccc|cccc|} \hline 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 0 \\ \hline G & C & G & C & C & G & G & A & A & G & G & A & T & G & A & C & G & G & A & G \\ \hline \end{array}$

- Divide  $B$  into chunks  $C_1, \dots, C_g$  of size  $\sigma$

# Data Structure for Batched Rank Queries

$$\mathbf{B} = \begin{array}{|cccc|cccc|cccc|cccc|cccc|} \hline 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 0 \\ \hline \mathbf{G} & \mathbf{C} & \mathbf{G} & \mathbf{C} & \mathbf{C} & \mathbf{G} & \mathbf{G} & \mathbf{A} & \mathbf{A} & \mathbf{G} & \mathbf{G} & \mathbf{A} & \mathbf{T} & \mathbf{G} & \mathbf{A} & \mathbf{C} & \mathbf{G} & \mathbf{G} & \mathbf{A} & \mathbf{G} \\ \hline \end{array}$$

$$M_A = 0010110101$$

$$M_C = 011010010$$

$$M_G = 011011011010111$$

$$M_T = 000010$$

- Divide  $\mathbf{B}$  into chunks  $C_1, \dots, C_g$  of size  $\sigma$
- For each character  $c$ , store  $M_c = 01^{d(1)} \dots 01^{d(g)}$  with select support
  - With  $d(i)$  the # of occurrences of  $c$  in  $C_i$

# Data Structure for Batched Rank Queries

$$B = \begin{array}{|cccc|cccc|cccc|cccc|cccc|} \hline 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 0 \\ \hline G & C & G & C & C & G & G & A & A & G & G & A & T & G & A & C & G & G & A & G \\ \hline \end{array}$$

$$M_A = 0010110101$$

$$M_C = 011010010$$

$$M_G = 011011011010111$$

$$M_T = 000010$$

$$R_1 = (C, 2, 1)(C, 4, 2)(G, 1, 1)(G, 3, 2)$$

$$R_2 = (A, 4, 1)(C, 1, 1)(G, 2, 1)(G, 3, 2)$$

$$R_3 = (A, 1, 1)(A, 4, 2)(G, 2, 1)(G, 3, 2)$$

$$R_4 = (A, 3, 1)(C, 4, 1)(G, 2, 1)(T, 1, 1)$$

$$R_5 = (A, 3, 1)(G, 1, 1)(G, 2, 1)(G, 4, 3)$$

- Divide  $B$  into chunks  $C_1, \dots, C_g$  of size  $\sigma$
- For each character  $c$ , store  $M_c = 01^{d(1)} \dots 01^{d(g)}$  with select support
  - With  $d(i)$  the # of occurrences of  $c$  in  $C_i$
- For each  $C_i[j] = c$ , store a triplet  $(c, j, \text{rank}_c(C, j))$  in  $R_i$ 
  - Sort triplets according by  $c$  and  $j$