

# Finding Optimal Alignment and Consensus of Circular Strings

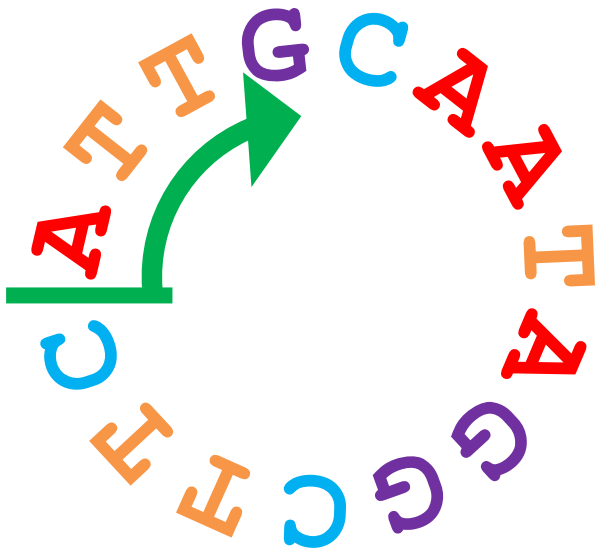
T. Lee\*, J.C. Na, H. Park, K. Park, and J.S. Sim

\*Seoul National University

June 23, 2010

# Circular String

- The first (leftmost) symbol is wrapped around next to the last (rightmost) symbol
- A circular string of length  $n$  can be read as  $n$  different *linear* strings (called *instances*)

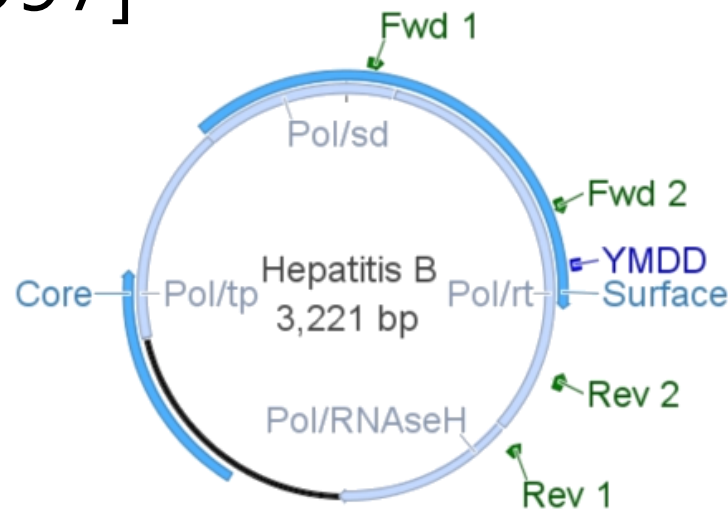


ATTGCAATAGGCTTC  
TTGCAATAGGCTTCA  
TGCAATAGGCTTCAT  
GCAATAGGCTTCA  
CAATAGGCTTCA

• • • • •

# Circular Strings in Nature

- "*Bacterial and mitochondrial DNA is typically circular, ... Consequently, tools for handling circular strings may someday be of use in those organisms.*"  
[Gusfield 1997]



HBV's genome sequence

# Consensus String Problem

- Given a set  $S$  of strings, find a representative string (consensus) of  $S$
- Application: Motif recognition, ...

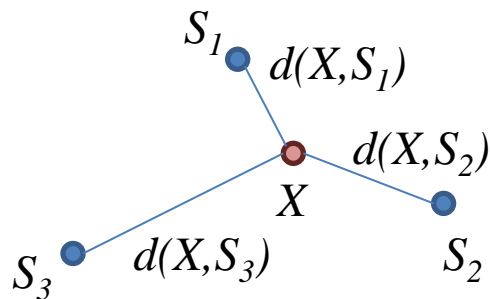
$S$  {  $s_1$  : a b c b a a c c a b e d a d a  
 $s_2$  : a a a b a b c c a b f d a c a  
 $s_3$  : a b a b d a c c a b e d a d a  
• • • • • • • • • • • • • •  
 $s_m$  : b b a b a a c a a b e d a d a

---

Consensus : a b a b a a c c a b e d a d a

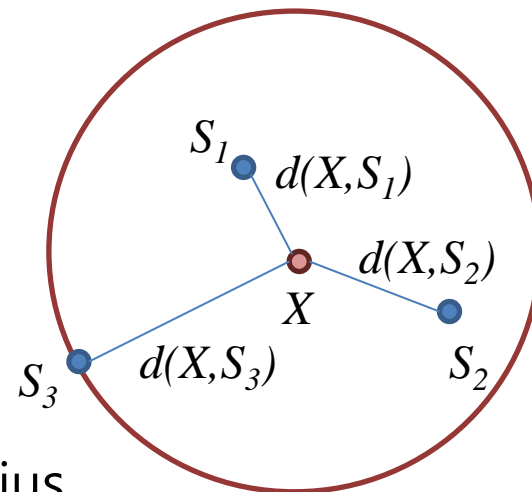
# Evaluation criteria

- Two major criteria of evaluating how well a string  $X$  represents a given set  $\mathcal{S}$  are
  - **Distance sum** : The sum of distances from string  $X$  to the strings in  $\mathcal{S}$
  - **Radius** : The longest distance from string  $X$  to the strings in  $\mathcal{S}$



distance sum

$$= d(X, S_1) + d(X, S_2) + d(X, S_3)$$



radius

$$= \max \{d(X, S_1), d(X, S_2), d(X, S_3)\}$$

# Evaluation criteria

- A good representative string should minimize distance sum **and/or** radius
  - Consensus minimizing distance sum (**CS**)
  - Consensus minimizing radius (**CR**)
  - Consensus minimizing both distance sum & radius (**CSR**)

# Distance measures

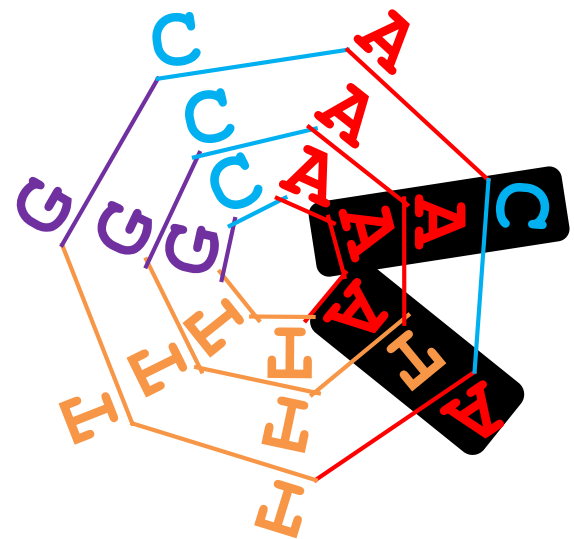
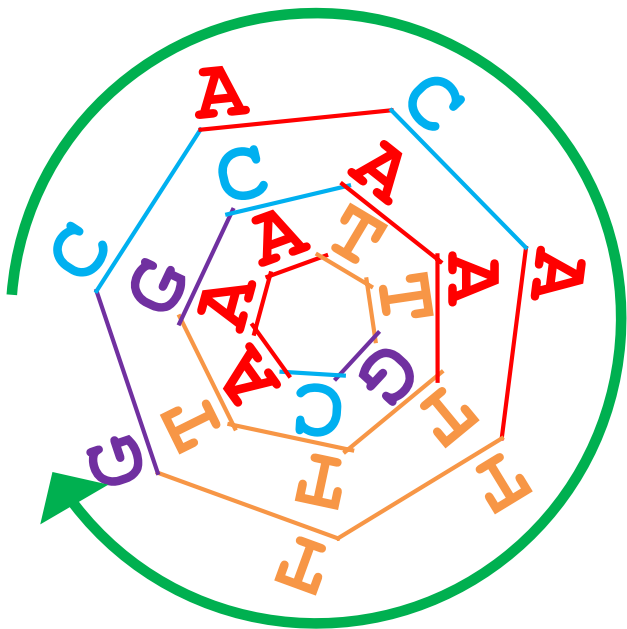
- Two well-known distances:
  - Hamming distance
    - Allowing only substitution
  - Edit distance
    - Allowing insertion and deletion as well as substitution

In this talk, we only consider the Hamming distance

# Consensus of Circular Strings

- Find not only **consensus** but also the optimal **alignment** of given set of circular strings

$S = \{CACATTG, GCAATTT, AATTGCA\}$



Opt. alignment w/  $ds=2$  &  $r=1$



# Previous Works

- Consensus of linear strings
  - CS: easy to find (in linear time)
    - Take majority symbol in each aligned position
  - CR: hard in general
    - NP-Complete for arbitrary number of strings (even for strings drawn from binary alphabet)
    - Approximation / Fixed parameter solutions

# Previous Works

- Fixed parameter solutions for CR and CSR
  - Algorithms for small constant  $m$  ( $=|S|$ )
  - [Gramm et al.](#) proposed a direct combinatorial algorithm for 3 strings
  - [Sze et al.](#) showed a condition for the existence of string whose radius is less than or equal to  $r$
  - [Boucher et al.](#) proposed an algorithm to find a string whose radius is  $\leq r$  for 4 binary strings.
  - [Amir et al.](#) proposed a linear time algorithm for CSR of 3 strings

# Previous Works

- Multiple alignment of circular strings  
[Mosig et al., Fernandes et al.]
  - The sum of pairs score
  - General purpose multiple sequence alignment techniques (e.g. clustalW)

# Our Contribution

- Goal
  - Find the optimal consensus and alignment of given set of circular strings
- We present efficient algorithms for 3 or 4 circular strings of length  $n$ 
  - $O(n^2 \log n)$  algorithm for CS of 3 strings
  - $O(n^2 \log n)$  algorithm for CR, CSR of 3 strings
  - $O(n^3 \log n)$  algorithm for CS of 4 strings

# Problem Definitions

- Circular string  $S$  of length  $n$ 
  - drawn from a constant-sized alphabet
  - The  $r$ th instance (linear string)  $S(r)$ 
    - $S(r) = S[r]S[(r+1) \bmod n] \dots S[(r-1) \bmod n]$
    - For circular string  $S = abcde$ ,  
 $S(0) = abcde$ ,  $S(1) = bcdea$ ,  $S(2) = cdeab$ , ...

# Problem Definitions

- $S = \{S_1, S_2, \dots, S_m\}$  of  $m$  circular strings of length  $n$
- Alignment  $\rho = (\rho_1, \rho_2, \dots, \rho_m)$  ( $\rho_k$ : integer in  $[0, n-1]$ )
  - Juxtaposition of  $S_1(\rho_1), S_2(\rho_2), \dots, S_m(\rho_m)$
  - At most  $n^{m-1}$  (not  $n^m$ ) distinct alignments
    - All  $(\rho_1 + k \bmod n, \rho_2 + k \bmod n, \dots, \rho_m + k \bmod n)$  are equivalent for  $k = 0, \dots, n-1$  under Hamming dist.

# Problem Definitions

- Given alignment  $\rho=(\rho_1, \rho_2, \dots, \rho_m)$  and a string  $X$ ,
  - Distance sum of  $X$  for  $\rho : E(\rho, X) = \sum_i d(X, S_i(\rho_i))$
  - Radius of  $X$  for  $\rho : R(\rho, X) = \max_i d(X, S_i(\rho_i))$
- For each alignment  $\rho$ ,
  - Min. distance sum for  $\rho : E_{\min}(\rho) = \min_{X'} E(\rho, X')$
  - Min. radius for  $\rho : R_{\min}(\rho) = \min_{X'} R(\rho, X')$
- For any alignment  $\rho'$  and string  $X'$ 
  - Optimal distance sum :  $E_{\text{opt}} = \min_{\rho', X'} E(\rho', X')$
  - Optimal radius :  $R_{\text{opt}} = \min_{\rho', X'} R(\rho', X')$

※  $d(x,y)$  : Hamming distance

# Problem Definitions

- Optimal consensus problems
  - Given a set  $\mathcal{S}$  of  $m$  circular strings of length  $n$ , find an optimal alignment  $\rho$  and a string  $X$  (*if any*) that satisfy:
    - Problem CS  $E(\rho, X) = E_{\text{opt}}$
    - Problem CR  $R(\rho, X) = R_{\text{opt}}$
    - Problem CSR  $E(\rho, X) = E_{\text{opt}}$  and  $R(\rho, X) = R_{\text{opt}}$

In this talk, we will use notations  $\text{CS}_k$ ,  $\text{CR}_k$ , and  $\text{CSR}_k$  to represent Problems CS, CR, and CSR of  $k$  circular strings



# Problem Definitions

- Bounded consensus problems
  - Given a set  $S$  of  $m$  circular strings of length  $n$ , and two integers  $s > 0$  and  $r > 0$ , find an ~~optimal~~ alignment  $\rho$  and a string  $X$  (if any) that satisfy:
    - Problem BS  $E(\rho, X) \leq s$
    - Problem BR  $R(\rho, X) \leq r$
    - Problem BSR  $E(\rho, X) \leq s$  and  $R(\rho, X) \leq r$

In this talk, we only present algorithms for CS and CSR, since they can be easily applied to CR, BS, BR, and BSR.

# Tools

- Algorithm for CSR of 3 linear strings  
[Amir et al. 09]
- Convolution method for counting matches/mismatches

# CSR of 3 Linear Strings

- For three linear strings of equal length  $n$ , CSR can be found in  $O(n)$  time

$s_1$	a	a	b	b	b	b	a	a	a	a	a	a	b	c
$s_2$	a	a	b	a	a	a	b	b	b	a	a	b	c	b
$s_3$	a	a	b	a	a	a	a	a	a	b	b	c	a	a

# CSR of 3 Linear Strings

- Every aligned position  $i$  is divided into five types and algorithm counts each Types  $k$  as  $c_k$

	Type 0	Type 1	Type 2	Type 3	Type 4									
$s_1$	a	a	b	b	b	a	a	a	a	a	a	a	b	c
$s_2$	a	a	b	a	a	a	b	b	b	a	a	b	c	b
$s_3$	a	a	b	a	a	a	a	a	a	b	b	c	a	a

$$c_0 = 3$$

$$c_1 = 3$$

$$c_2 = 3$$

$$c_3 = 2$$

$$c_4 = 3$$

- Type 0:  $s_1[i] = s_2[i] = s_3[i]$  (all matches)
- Type 1:  $s_1[i] \neq s_2[i] = s_3[i]$  ( $s_1[i]$  is the minority)
- Type 2:  $s_2[i] \neq s_1[i] = s_3[i]$  ( $s_2[i]$  is the minority)
- Type 3:  $s_3[i] \neq s_1[i] = s_2[i]$  ( $s_3[i]$  is the minority)
- Type 4:  $s_1[i] \neq s_2[i] \neq s_3[i]$  (all mismatches)

# CSR of 3 Linear Strings

- From  $c_1 \sim c_4$ , the minimum possible distance sum  $E_{min}$  and radius  $R_{min}$  can be determined:
  - $E_{min} = c_1 + c_2 + c_3 + 2c_4$
  - $R_{min} \geq \max(L_1, L_2)$  ( $L_1 = \lceil \max_{i \neq j} d(s_i, s_j) / 2 \rceil$ ,  $L_2 = \lceil E_{min} / 3 \rceil$ )

	Type 0			Type 1			Type 2			Type 3		Type 4		
$s_1$	a	a	b	b	b	b	a	a	a	a	a	a	b	c
$s_2$	a	a	b	a	a	a	b	b	b	a	a	b	c	b
$s_3$	a	a	b	a	a	a	a	a	a	b	b	c	a	a
$s$	a	a	b	a	a	a	a	a	a	a	a	?	?	?

$$L_1 = 5, \quad L_2 = 5$$

$$E_{min} = 14, \quad R_{min} = 5$$

# CSR of 3 Linear Strings

- Algorithm then finds (constructs) a consensus string  $s$  with  $E_{min}$  and radius  $R_{min}$  if it exists
  - By selecting Type 4 symbols wisely to balance the distances of strings from the consensus.

	Type 0	Type 1	Type 2	Type 3	Type 4
$s_1$	a a b	b b b	a a a	a a	a b c
$s_2$	a a b	a a a	b b b	a a	b c b
$s_3$	a a b	a a a	a a a	b b	c a a
$s$	a a b	a a a	a a a	a a a	b c c

$$E_{min} = 14, \quad R_{min} = 5_2$$

# Tools

- **Lemma 1.** For three linear strings of length  $n$ , if we are given counters  $c_1$  to  $c_4$ ,
  - $E_{min}$ ,  $R_{min}$  and the existence of a consensus w/ both  $E_{min}$  and  $R_{min}$  are determined in  $O(1)$  time
  - We can construct such a consensus in  $O(n)$  time
- Lemma 2. Given two circular strings  $X$  and  $Y$  of equal length  $n$  ( $|\Sigma|$  is constant), the numbers of matches or mismatches in all  $n$  alignments can be computed in  $O(n \log n)$  time, even with the presence of don't care or mismatch symbols.

# Discrete Convolution

Given integer arrays  $t$  and  $p$  ( $|t|=n$  and  $|p|=m$ ),  $t \otimes p$  is the array of all inner products, where:

$$(t \otimes p)[i] = \sum_{j=0..m-1} t[i+j] p[j] \quad (i=0, \dots, n-m).$$

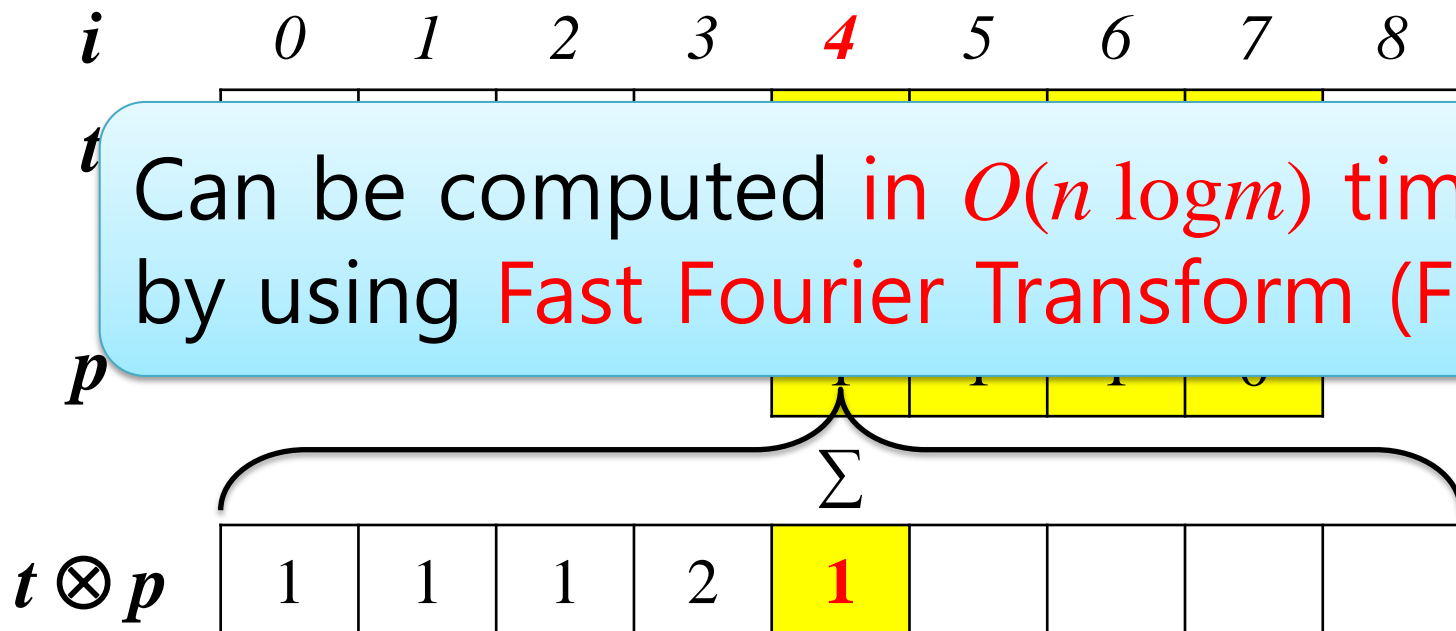
$i$	$0$	$1$	$2$	<b>3</b>	$4$	$5$	$6$	$7$	$8$
$t$	1	0	0	1	0	1	0	0	0
				×	×	×	×		
$p$				1	1	1	0		
				Σ					
$t \otimes p$	1	1	1	<b>2</b>					



# Discrete Convolution

Given integer arrays  $t$  and  $p$  ( $|t|=n$  and  $|p|=m$ ),  $t \otimes p$  is the array of all inner products, where:

$$(t \otimes p)[i] = \sum_{j=0..m-1} t[i+j] p[j] \quad (i=0, \dots, n-m).$$



# Counting Matches/Mismatches

- Using discrete convolution, we can count matches/mismatches in **every** alignments of two circular strings efficiently
  - in  **$O(n \log n)$**  time ( $|\Sigma|$  is constant)

Alignment	$(0, 0)$	$(0, 1)$	$(0, 2)$	$(0, 3)$
X	a a b a	a a b a	a a b a	a a b a
Y	a b b a	b b a a	b a a b	a a b b
<b>Matches</b>	<b>3</b>	<b>1</b>	<b>1</b>	<b>3</b>
<b>Mismatches</b>	<b>1</b>	<b>3</b>	<b>3</b>	<b>1</b>

# Counting Matches/Mismatches

- Bit mask  $B_{x,\sigma}$  and inverse bit mask  $B'_{x,\sigma}$ 
  - For given string  $x$  and each symbol  $\sigma \in \Sigma$ ,

$$B_{x,\sigma}[i] = \begin{cases} 1, & \text{if } x[i] = \sigma \\ 0, & \text{if } x[i] \neq \sigma \end{cases} \quad \text{and} \quad B'_{x,\sigma}[i] = \begin{cases} 0, & \text{if } x[i] = \sigma \\ 1, & \text{if } x[i] \neq \sigma \end{cases}$$

$x$	<b>a</b>	<b>a</b>	<b>b</b>	<b>a</b>
$B_{x,\mathbf{a}}$	<b>1</b>	<b>1</b>	<b>0</b>	<b>1</b>
$B_{x,\mathbf{b}}$	<b>0</b>	<b>0</b>	<b>1</b>	<b>0</b>

# Counting Matches/Mismatches

	$i$	0	1	2	3	4	5	6	7
$X(0)X(0)$ $\rightarrow$ $X^2(0)$		a	a	b	a	a	a	b	a
$Y(0)$		a	b	b	a				
$B_{X^2(0),a}$		1	1	0	1	1	1	0	1
$B_{Y(0),a}$		1	0	0	1				
$B_{X^2(0),a} \otimes B_{Y(0),a}$		2							
$B_{X^2(0),b}$		0	0	1	0	0	0	1	0
$B_{Y(0),b}$		0	1	1	0				
$B_{X^2(0),b} \otimes B_{Y(0),b}$		1							
$\sum_{\sigma} B_{X^2(0),\sigma} \otimes B_{Y(0),\sigma}$		3							

#. matches in alignment (0,0)

# Counting Matches/Mismatches

$i$	0	1	2	3	4	5	6	7
$X^2(0)$	a	a	b	a	a	a	b	a
$Y(0)$		a	b	b	a			
$B_{X^2(0),a}$	1	1	0	1	1	1	0	1
$B_{Y(0),a}$		1	0	0	1			
$B_{X^2(0),a} \otimes B_{Y(0),a}$	2	2						
$B_{X^2(0),b}$	0	0	1	0	0	0	1	0
$B_{Y(0),b}$		0	1	1	0			
$B_{X^2(0),b} \otimes B_{Y(0),b}$	1	1						
$\sum_{\sigma} B_{X^2(0),\sigma} \otimes B_{Y(0),\sigma}$	3	3						

#. matches in alignment (0,3)

# Counting Matches/Mismatches

$i$	0	1	2	3	4	5	6	7
$X^2(0)$	a	a	b	a	a	a	b	a
$Y(0)$			a	b	b	a		
$B_{X^2(0),a}$	1	1	0	1	1	1	0	1
$B_{Y(0),a}$			1	0	0	1		
$B_{X^2(0),a} \otimes B_{Y(0),a}$	2	2	1					
$B_{X^2(0),b}$	0	0	1	0	0	0	1	0
$B_{Y(0),b}$			0	1	1	0		
$B_{X^2(0),b} \otimes B_{Y(0),b}$	1	1	0					
$\sum_{\sigma} B_{X^2(0),\sigma} \otimes B_{Y(0),\sigma}$	3	3	1					

#. matches in alignment (0,2)

# Counting Matches/Mismatches

$i$	0	1	2	3	4	5	6	7
$X^2(0)$	a	a	b	a	a	a	b	a
$Y(0)$				a	b	b	a	
$B_{X^2(0),a}$	1	1	0	1	1	1	0	1
$B_{Y(0),a}$				1	0	0	1	
$B_{X^2(0),a} \otimes B_{Y(0),a}$	2	2	1	1				
$B_{X^2(0),b}$	0	0	1	0	0	0	1	0
$B_{Y(0),b}$				0	1	1	0	
$B_{X^2(0),b} \otimes B_{Y(0),b}$	1	1	0	0				
$\sum_{\sigma} B_{X^2(0),\sigma} \otimes B_{Y(0),\sigma}$	3	3	1	1				

#. matches in alignment (0,1)

# Counting Matches/Mismatches

- Counting matches with mismatch '#'
  - $B_{x,\sigma}[i]=0$  if  $x[i]=\#$  and  $B_{y,\sigma}[i]=0$  if  $y[i]=\#$  for  $\forall \sigma$
- Counting mismatches with don't care '\$'
  - $B_{x,\sigma}[i]=0$  if  $x[i]=\$$  and  $B'_{y,\sigma}[i]=0$  if  $y[i]=\$$  for  $\forall \sigma$



# Tools

- Lemma 1. For three linear strings of length  $n$ , if we are given counters  $c_1$  to  $c_4$ ,
  - $E_{min}$ ,  $R_{min}$  and the existence of a consensus w/ both  $E_{min}$  and  $R_{min}$  are determined in  $O(1)$  time
  - We can construct such a consensus in  $O(n)$  time
- **Lemma 2.** Given two circular strings  $X$  and  $Y$  of equal length  $n$  ( $|\Sigma|$  is constant), the numbers of matches or mismatches in all  $n$  alignments can be computed **in  $O(n \log n)$  time**, even with the presence of don't care or mismatch symbols.

# Algorithms

- A naïve algorithm
  - For CS3 or CSR3, we apply  $O(n)$ -time algorithm for linear strings in each of all  $O(n^2)$  alignments  $\rightarrow O(n^3)$  time
  - Similarly, CS4 can be found in  $O(n^4)$  time
- Our algorithms achieve  $(n/\log n)$ -speedup
  - $O(n^2 \log n)$  time for CS3 and CSR3
  - $O(n^3 \log n)$  time for CS4

# Algorithm for Problem CS3

- Computation of  $E_{\min}(0, \delta, \gamma)$ 
  - We first superpose  $S_1(0)$  and  $S_2(\delta)$  into  $Z_\delta$ 
$$Z_\delta[i] = \langle S_1(0)[i], S_2(\delta)[i] \rangle$$
  - Generate bit masks of  $Z_\delta$  and  $S_3(\gamma)$ 
$$B_{Z_\delta\sigma}[i] = 1, \text{ if } Z_\delta[i] = \langle \sigma, * \rangle \text{ or } \langle *, \sigma \rangle;$$
$$0, \text{ otherwise}$$
  - Compute inner products of  $B_{Z_\delta\sigma}$  and  $B'_{S_3(\gamma),\sigma}$
  - Key observation: sum of all inner products over  $\sigma$  equals to  $E_{\min}(0, \delta, \gamma)$

# Algorithm for Problem CS3

- For each  $\sigma$ ,  $B_{(Z_\delta)^2, \sigma} \otimes B'_{S_{\mathcal{Z}(0)}, \sigma}$  produces the array of  $n$  inner products  $B_{(Z_\delta)^2, \sigma}$  and  $B'_{S_{\mathcal{Z}(\gamma)}, \sigma}$  ( $\gamma=0, \dots, n-1$ )
- Thus, we can compute  $n$  distance sums  $E_{\min}(0, \delta, 0) \sim E_{\min}(0, \delta, n-1)$  in  $O(n \log n)$  time by FFT

Row 0	$E_{\min}(0, 0, 0)$	$E_{\min}(0, 0, 1)$	***	$E_{\min}(0, 0, n-1)$
***			***	
Row $\delta$	$E_{\min}(0, \delta, 0)$	$E_{\min}(0, \delta, 1)$	***	$E_{\min}(0, \delta, n-1)$
***			***	
Row $n-1$	$E_{\min}(0, n-1, 0)$	$E_{\min}(0, n-1, 1)$	***	$E_{\min}(0, n-1, n-1)$

$E_{\min}$ 's for all  $n^2$  alignment  $\rightarrow O(n^2 \log n)$  time

# Algorithm for Problem CS3

- Workflow
  - Compute  $E_{\min}$  for all alignments in  $O(n^2 \log n)$  time by using FFT
  - Among all  $n^2$  alignments, find the best alignment  $\rho$  with the smallest  $E_{\min}$  ( $=E_{\text{opt}}$ ) and construct a consensus in  $\rho$
- **Theorem 1.** Problem CS3 can be solved in  $O(n^2 \log n)$  time and  $O(n)$  space

# Algorithm for Problem CSR3

- Algorithm for CS3 directly computes  $E_{\min}$  for each alignment, but it does not give any information of **counters  $c_1$  to  $c_4$**  which are essential to apply Lemma 1
- How can we efficiently count  $c_1$  to  $c_4$  for each alignment?
  - Convolution and system of linear equations

# Algorithm for Problem CSR3

- For three linear strings, we obtain :

$$\left\{ \begin{array}{l} c_1 + c_2 + c_3 + c_4 = n - c_0, \\ c_1 + c_2 + c_4 = d(s_1, s_2), \\ c_1 + c_3 + c_4 = d(s_1, s_3), \\ c_2 + c_3 + c_4 = d(s_2, s_3). \end{array} \right.$$

	Type 0	Type 1	Type 2	Type 3	Type 4
$s_1$	a a b	b b b	a a a	a a	a b c
$s_2$	a a b	a a a	b b b	a a	b c b
$s_3$	a a b	a a a	a a a	b b	c a a
$d(s_1, s_2) =$		3 +	3 +		3

# Algorithm for Problem CSR3

- Once  $c_0$  and the **pairwise distances** have been computed, we can compute  $c_1$  to  $c_4$

$$c_1 = n - d(s_2, s_3) - c_0,$$

$$c_2 = n - d(s_1, s_3) - c_0,$$

$$c_3 = n - d(s_1, s_2) - c_0,$$

$$c_4 = d(s_1, s_2) + d(s_1, s_3) + d(s_2, s_3) - 2c_0 - 2n$$































































# Algorithm for Problem CSR3

- Algorithm efficiently computes
  - $d(S_i(0), S_j(\delta))$ 's for all  $i, j, \delta$  in  $O(n \log n)$  time using FFT (by Lemma 2)
  - Counters  $c_0$ 's for all alignments in  $O(n^2 \log n)$  time using FFT and superposition (similar to Algorithm CS3)
- **Theorem 2.** Problem CSR3 can be solved in  $O(n^2 \log n)$  time and  $O(n)$  space

# Algorithm for Problem CS4

- Given four linear strings (instances),
  - We define 5 types and 15 counters wrt. combinations of distinct symbols
  - $E_{\min} = \sum_i b_i + 2(\sum_i c_i + \sum_i d_i) + 3e$

Type	A	B				C			D						E
$s_1$															
$s_2$															
$s_3$															
$s_4$															
Counter	$a$	$b_1$	$b_2$	$b_3$	$b_4$	$c_1$	$c_2$	$c_3$	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$	$e$

# Algorithm for Problem CS4

- System of linear equations

$$a + b_3 + b_4 + c_3 + d_1 = M_{12}, \quad a + b_4 = M_{123}, \quad a = M_{1234},$$

$$a + b_2 + b_4 + c_2 + d_2 = M_{13}, \quad a + b_3 = M_{124},$$

$$a + b_2 + b_3 + c_1 + d_3 = M_{14}, \quad a + b_2 = M_{134},$$

$$a + b_1 + b_4 + c_1 + d_4 = M_{23}, \quad a + b_1 = M_{234},$$

$$a + b_1 + b_3 + c_2 + d_5 = M_{24},$$

$$a + b_1 + b_2 + c_3 + d_6 = M_{34},$$

$$c_3 = d(s_{12}, s_{34}), \quad c_2 = d(s_{13}, s_{24}), \quad c_1 = d(s_{14}, s_{23}),$$

$$a + \sum b_i + \sum c_i + \sum d_i + e = n .$$

# Algorithm for Problem CS4

- $M_{ij}$ ,  $M_{ijk}$ ,  $M_{ijkl}$ : Numbers of pair/triple/quadruple-wise matches ( $\{i,j,k,l\}=\{1,2,3,4\}$ )
- Pair string  $s_{ij}$  of  $s_i$  and  $s_j$ 
  - $s_{ij}[p] = s_i[p]$  if  $s_i[p] = s_j[p]$ ; otherwise,  $s_{ij}[p] = \$$
- Distance between  $s_{ij}$  and  $s_{kl}$ ,  $d(s_{ij}, s_{kl})$

$s_1$	a	b	a	a	a	a	a	b	b	} $M_{12} = 5$
$s_2$	a	c	a	a	b	b	b	b	b	
$s_3$	b	a	a	a	a	c	b	a	c	
$s_4$	b	a	b	a	b	d	a	b	c	
$s_{12}$	a	\$	a	a	\$	\$	\$	b	b	}
$s_{34}$	b	a	\$	a	\$	\$	\$	\$	c	

# Algorithm for Problem CS4

- $M_{ij}$ ,  $M_{ijk}$ ,  $M_{ijkl}$ : Numbers of pair/triple/quadruple-wise matches ( $\{i,j,k,l\}=\{1,2,3,4\}$ )
- Pair string  $s_{ij}$  of  $s_i$  and  $s_j$ 
  - $s_{ij}[p] = s_i[p]$  if  $s_i[p] = s_j[p]$ ; otherwise,  $s_{ij}[p] = \$$
- Distance between  $s_{ij}$  and  $s_{kl}$ ,  $d(s_{ij}, s_{kl})$

$s_1$	a	b	a	a	a	a	a	b	b	} $M_{12} = 5$	
$s_2$	a	c	a	a	b	b	b	b	b		} $M_{124} = 2$
$s_3$	b	a	a	a	a	c	b	a	c		
$s_4$	b	a	b	a	b	d	a	b	c		
$s_{12}$	a	\$	a	a	\$	\$	\$	b	b	} $M_{124} = 2$	
$s_{34}$	b	a	\$	a	\$	\$	\$	\$	c		

# Algorithm for Problem CS4

- $M_{ij}$ ,  $M_{ijk}$ ,  $M_{ijkl}$ : Numbers of pair/triple/quadruple-wise matches ( $\{i,j,k,l\}=\{1,2,3,4\}$ )
- Pair string  $s_{ij}$  of  $s_i$  and  $s_j$ 
  - $s_{ij}[p] = s_i[p]$  if  $s_i[p] = s_j[p]$ ; otherwise,  $s_{ij}[p] = \$$
- Distance between  $s_{ij}$  and  $s_{kl}$ ,  $d(s_{ij}, s_{kl})$

$s_1$	a	b	a	a	a	a	a	b	b	} $M_{12} = 5$		
$s_2$	a	c	a	a	b	b	b	b	b		} $M_{124} = 2$	
$s_3$	b	a	a	a	a	c	b	a	c			} $M_{1234} = 1$
$s_4$	b	a	b	a	b	d	a	b	c			
$s_{12}$	a	\$	a	a	\$	\$	\$	b	b	} $M_{1234} = 1$		
$s_{34}$	b	a	\$	a	\$	\$	\$	\$	c			

# Algorithm for Problem CS4

- $M_{ij}$ ,  $M_{ijk}$ ,  $M_{ijkl}$ : Numbers of pair/triple/quadruple-wise matches ( $\{i,j,k,l\}=\{1,2,3,4\}$ )
- Pair string  $s_{ij}$  of  $s_i$  and  $s_j$ 
  - $s_{ij}[p] = s_i[p]$  if  $s_i[p] = s_j[p]$ ; otherwise,  $s_{ij}[p] = \$$
- Distance between  $s_{ij}$  and  $s_{kl}$ ,  $d(s_{ij}, s_{kl})$

$s_1$	a	b	a	a	a	a	a	b	b	} $M_{12} = 5$	
$s_2$	a	c	a	a	b	b	b	b	b		} $M_{124} = 2$
$s_3$	b	a	a	a	a	c	b	a	c		
$s_4$	b	a	b	a	b	d	a	b	c		
$s_{12}$	a	\$	a	a	\$	\$	\$	b	b	} $d(s_{12}, s_{34}) = 2$	
$s_{34}$	b	a	\$	a	\$	\$	\$	\$	c		

# Algorithm for Problem CS4

- By using FFT, we can efficiently compute, for all  $n^3$  alignments,
  - all  $M_{ij}$ ,  $M_{ijk}$ ,  $M_{ijkl}$  in  $O(n^2)$ ,  $O(n^3)$ , and  $O(n^3 \log n)$  time
  - all  $d(s_{ij}, s_{kl})$  in  $O(n^3 \log n)$  time
- **Theorem 3.** Problem CS4 can be solved in  $O(n^3 \log n)$  time and  $O(n)$  space



# Conclusion

- We proposed efficient algorithms to solve the optimal/bounded consensus problems for 3 or 4 circular strings
- Future works
  - Algorithms for CSR4 or more strings
  - Algorithms with edit distance measure

**Thank you!**