

# Consensus Strings with Small Maximum Distance and Small Distance Sum

Laurent Bulteau<sup>1</sup>, **Markus L. Schmid**<sup>2</sup>

<sup>1</sup> Université Paris-Est, France

<sup>2</sup> Trier University, Germany

MFCS 2018

## Consensus String Problems

Input: A set (multi-set) of strings.

Output: A string that is a good **consensus** of the input strings.

# Consensus String Problems

Input: A set (multi-set) of strings.

Output: A string that is a good **consensus** of the input strings.

$s_1$	c	b	c	a	b	a	a
$s_2$	c	b	c	a	b	c	b
$s_3$	a	b	c	c	c	c	a
$s_4$	c	c	c	a	b	c	a
$s_5$	c	b	c	a	a	c	a
$s_6$	c	b	c	a	a	c	a
$s_7$	a	b	b	a	b	a	a
$s_8$	b	b	c	a	a	c	a

---

# Consensus String Problems

Input: A set (multi-set) of strings.

Output: A string that is a good **consensus** of the input strings.

$s_1$     c b c a b a a

$s_2$     c b c a b c b

$s_3$     a b c c c c a

$s_4$     c c c a b c a

$s_5$     c b c a a c a

$s_6$     c b c a a c a

$s_7$     a b b a b a a

$s_8$     b b c a a c a

---

$s$     c b c a b c a

# Consensus String Problems

Input: A set (multi-set) of strings.

Output: A string that is a good **consensus** of the input strings.

$s_1$	c	b	c	a	b	a	a
$s_2$	c	b	c	a	b	c	b
$s_3$	a	b	c	c	c	c	a
$s_4$	c	c	c	a	b	c	a
$s_5$	c	b	c	a	a	c	a
$s_6$	c	b	c	a	a	c	a
$s_7$	a	b	b	a	b	a	a
$s_8$	b	b	c	a	a	c	a
<hr/>							
$s$	c	b	c	a	b	c	a

# Consensus String Problems

Input: A set (multi-set) of strings.

Output: A string that is a good **consensus** of the input strings.

$s_1$	c	b	c	a	b	a	a
$s_2$	c	b	c	a	b	c	b
$s_3$	a	b	c	c	c	c	a
$s_4$	c	c	c	a	b	c	a
$s_5$	c	b	c	a	a	c	a
$s_6$	c	b	c	a	a	c	a
$s_7$	a	b	b	a	b	a	a
$s_8$	b	b	c	a	a	c	a
<hr/>							
$s$	a	b	c	a	b	c	a

# Basic Notations

Standard string notations:

$\Sigma$  finite alphabet

$\Sigma^*$  words over  $\Sigma$

$\Sigma^n$   $\{w \in \Sigma^* \mid |w| = n\}$

$\Sigma^{\leq n}$   $\bigcup_{i=0}^n \Sigma^i$

$d_H(u, v)$  Hamming distance

$\preceq$  substring relation,

$u \preceq v \Leftrightarrow v = xuy$

## Basic Notations

Standard string notations:

$\Sigma$  finite alphabet

$\Sigma^*$  words over  $\Sigma$

$\Sigma^n$   $\{w \in \Sigma^* \mid |w| = n\}$

$\Sigma^{\leq n}$   $\bigcup_{i=0}^n \Sigma^i$

$d_H(u, v)$  Hamming distance

$\preceq$  substring relation,

$u \preceq v \Leftrightarrow v = xuy$

For multi-set  $S \subseteq \Sigma^\ell$  and  $v \in \Sigma^\ell$ :

$r_H(v, S) = \max\{d_H(v, u) \mid u \in S\}$  **radius** of  $S$  (w. r. t.  $v$ )

$s_H(v, S) = \sum_{u \in S} d_H(v, u)$  **distance sum** of  $S$  (w. r. t.  $v$ )



# The Closest String Problem

## $(r, s)$ -CLOSEST STRING

*Instance:* Multi-set  $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^\ell$ ,  $\ell \in \mathbb{N}$ ,  
 $d_r, d_s \in \mathbb{N}$ .

*Question:* Is there an  $s \in \Sigma^\ell$  with  
 $r_H(s, S) \leq d_r$  and  $s_H(s, S) \leq d_s$ ?

# The Closest String Problem

## $(r, s)$ -CLOSEST STRING

*Instance:* Multi-set  $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^\ell$ ,  $\ell \in \mathbb{N}$ ,  
 $d_r, d_s \in \mathbb{N}$ .

*Question:* Is there an  $s \in \Sigma^\ell$  with  
 $r_H(s, S) \leq d_r$  and  $s_H(s, S) \leq d_s$ ?

$s_1$	c	b	c	a	b	a	a
$s_2$	c	b	c	a	b	c	b
$s_3$	a	b	c	c	c	c	a
$s_4$	c	c	c	a	b	c	a
$s_5$	c	b	c	a	a	c	a
$s_6$	c	b	c	a	a	c	a
$s_7$	a	b	b	a	b	a	a
$s_8$	b	b	c	a	a	c	a

---

$k = 8$

$\ell = 7$

$d_r = 2$

$d_s = 20$

# The Closest String Problem

## (r, s)-CLOSEST STRING

*Instance:* Multi-set  $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^\ell$ ,  $\ell \in \mathbb{N}$ ,  
 $d_r, d_s \in \mathbb{N}$ .

*Question:* Is there an  $s \in \Sigma^\ell$  with  
 $r_H(s, S) \leq d_r$  and  $s_H(s, S) \leq d_s$ ?

$s_1$    c b c a b a a

$s_2$    c b c a b c b

$s_3$    a b c c c c a

$s_4$    c c c a b c a

$s_5$    c b c a a c a

$s_6$    c b c a a c a

$s_7$    a b b a b a a

$s_8$    b b c a a c a

---

$s$    a b c a b c a

$k = 8$

$\ell = 7$

$d_r = 2$

$d_s = 20$

# The Closest String Problem

## $(r, s)$ -CLOSEST STRING

*Instance:* Multi-set  $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^\ell$ ,  $\ell \in \mathbb{N}$ ,  
 $d_r, d_s \in \mathbb{N}$ .

*Question:* Is there an  $s \in \Sigma^\ell$  with  
 $r_H(s, S) \leq d_r$  and  $s_H(s, S) \leq d_s$ ?

$s_1$	c	b	c	a	b	a	a
$s_2$	c	b	c	a	b	c	b
$s_3$	a	b	c	c	c	c	a
$s_4$	c	c	c	a	b	c	a
$s_5$	c	b	c	a	a	c	a
$s_6$	c	b	c	a	a	c	a
$s_7$	a	b	b	a	b	a	a
$s_8$	b	b	c	a	a	c	a
<hr/>							
$s$	a	b	c	a	b	c	a

$$k = 8$$

$$\ell = 7$$

$$d_r = 2$$

$$d_s = 20$$

$$r_H(s, S) = 2$$

$$s_H(s, S) = 16$$

# The “Substring Variant”

## (r, s)-CLOSEST SUBSTRING

*Instance:* Multi-set  $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^{\leq \ell}$ ,  $\ell \in \mathbb{N}$ ,  
 $d_r, d_s, m \in \mathbb{N}$ .

*Question:* Is there an  $s \in \Sigma^m$ ,  
multi-set  $S' = \{s'_i \mid s'_i \preceq s_i, 1 \leq i \leq k\} \subseteq \Sigma^m$  with  
 $r_H(s, S') \leq d_r$  and  $s_H(s, S') \leq d_s$ ?

# The “Substring Variant”

## (r, s)-CLOSEST SUBSTRING

*Instance:* Multi-set  $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^{\leq \ell}$ ,  $\ell \in \mathbb{N}$ ,  
 $d_r, d_s, m \in \mathbb{N}$ .

*Question:* Is there an  $s \in \Sigma^m$ ,  
multi-set  $S' = \{s'_i \mid s'_i \preceq s_i, 1 \leq i \leq k\} \subseteq \Sigma^m$  with  
 $r_H(s, S') \leq d_r$  and  $s_H(s, S') \leq d_s$ ?

$s_1$  a a c b c a b a a

$s_2$  b c b c a b c b

$s_3$  a a b c c

$s_4$  c c c a b c a c

$s_5$  c c b c a a c a

$s_6$  a c b c a a

$s_7$  a a b b a b a a

$s_8$  b b b c a a c a c c b

---

$k = 8$

$\ell = 11$

$m = 4$

$d_r = 3$

$d_s = 7$

# The “Substring Variant”

## (r, s)-CLOSEST SUBSTRING

*Instance:* Multi-set  $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^{\leq \ell}$ ,  $\ell \in \mathbb{N}$ ,  
 $d_r, d_s, m \in \mathbb{N}$ .

*Question:* Is there an  $s \in \Sigma^m$ ,  
multi-set  $S' = \{s'_i \mid s'_i \preceq s_i, 1 \leq i \leq k\} \subseteq \Sigma^m$  with  
 $r_H(s, S') \leq d_r$  and  $s_H(s, S') \leq d_s$ ?

$s_1$  a a c b c a b a a

$s_2$  b c b c a b c b

$s_3$  a a b c c

$s_4$  c c c a b c a c

$s_5$  c c b c a a c a

$s_6$  a c b c a a

$s_7$  a a b b a b a a

$s_8$  b b b c a a c a c c b

---

$s$  a b c a

$k = 8$

$\ell = 11$

$m = 4$

$d_r = 3$

$d_s = 7$

# The “Substring Variant”

## (r, s)-CLOSEST SUBSTRING

*Instance:* Multi-set  $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^{\leq \ell}$ ,  $\ell \in \mathbb{N}$ ,  
 $d_r, d_s, m \in \mathbb{N}$ .

*Question:* Is there an  $s \in \Sigma^m$ ,  
multi-set  $S' = \{s'_i \mid s'_i \preceq s_i, 1 \leq i \leq k\} \subseteq \Sigma^m$  with  
 $r_H(s, S') \leq d_r$  and  $s_H(s, S') \leq d_s$ ?

$s_1$  a a c b c a b a a

$s_2$  b c b c a b c b

$s_3$  a a b c c

$s_4$  c c c a b c a c

$s_5$  c c b c a a c a

$s_6$  a c b c a a

$s_7$  a a b b a b a a

$s_8$  b b b c a a c a c c b

---

$s$  a b c a

$k = 8$

$\ell = 11$

$m = 4$

$d_r = 3$

$d_s = 7$



# The “Substring Variant”

## (r, s)-CLOSEST SUBSTRING

*Instance:* Multi-set  $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^{\leq \ell}$ ,  $\ell \in \mathbb{N}$ ,  
 $d_r, d_s, m \in \mathbb{N}$ .

*Question:* Is there an  $s \in \Sigma^m$ ,  
multi-set  $S' = \{s'_i \mid s'_i \preceq s_i, 1 \leq i \leq k\} \subseteq \Sigma^m$  with  
 $r_H(s, S') \leq d_r$  and  $s_H(s, S') \leq d_s$ ?

$s_1$  a a c b c a b a a

$s_2$  b c b c a b c b

$s_3$  a a b c c

$s_4$  c c c a b c a c

$s_5$  c c b c a a c a

$s_6$  a c b c a a

$s_7$  a a b b a b a a

$s_8$  b b b c a a c a c c b

---

$s$  a b c a

$k = 8$

$\ell = 11$

$m = 4$

$d_r = 3$

$d_s = 7$

$r_H(s, S') = 1$

$s_H(s, S') = 7$

# The “Outlier Variant”

## (r, s)-CLOSEST STRING WITH OUTLIERS

*Instance:* Multi-set  $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^\ell$ ,  $\ell \in \mathbb{N}$ ,  
 $d_r, d_s, t \in \mathbb{N}$ .

*Question:* Is there an  $s \in \Sigma^\ell$ ,  
 $S' \subseteq S$  with  $|S'| = k - t$  with  
 $r_H(s, S') \leq d_r$  and  $s_H(s, S') \leq d_s$ ?

# The “Outlier Variant”

## (r, s)-CLOSEST STRING WITH OUTLIERS

*Instance:* Multi-set  $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^\ell$ ,  $\ell \in \mathbb{N}$ ,  
 $d_r, d_s, t \in \mathbb{N}$ .

*Question:* Is there an  $s \in \Sigma^\ell$ ,  
 $S' \subseteq S$  with  $|S'| = k - t$  with  
 $r_H(s, S') \leq d_r$  and  $s_H(s, S') \leq d_s$ ?

$s_1$     c   b   c   a   b   a   a

$s_2$     c   b   c   a   b   c   b

$s_3$     a   b   c   c   c   c   a

$s_4$     c   c   c   a   b   c   a

$s_5$     c   b   c   a   a   c   a

$s_6$     c   b   c   a   a   c   a

$s_7$     a   b   b   a   b   a   a

$s_8$     b   b   c   a   a   c   a

---

$k = 8$

$\ell = 7$

$t = 2$

$d_r = 2$

$d_s = 8$

# The “Outlier Variant”

## (r, s)-CLOSEST STRING WITH OUTLIERS

*Instance:* Multi-set  $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^\ell$ ,  $\ell \in \mathbb{N}$ ,  
 $d_r, d_s, t \in \mathbb{N}$ .

*Question:* Is there an  $s \in \Sigma^\ell$ ,  
 $S' \subseteq S$  with  $|S'| = k - t$  with  
 $r_H(s, S') \leq d_r$  and  $s_H(s, S') \leq d_s$ ?

$s_1$     c b c a b a a

$s_2$     c b c a b c b

$s_4$     c c c a b c a

$s_5$     c b c a a c a

$s_6$     c b c a a c a

$s_8$     b b c a a c a

$k = 8$

$\ell = 7$

$t = 2$

$d_r = 2$

$d_s = 8$

# The “Outlier Variant”


## (r, s)-CLOSEST STRING WITH OUTLIERS

*Instance:* Multi-set  $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^\ell$ ,  $\ell \in \mathbb{N}$ ,  
 $d_r, d_s, t \in \mathbb{N}$ .

*Question:* Is there an  $s \in \Sigma^\ell$ ,  
 $S' \subseteq S$  with  $|S'| = k - t$  with  
 $r_H(s, S') \leq d_r$  and  $s_H(s, S') \leq d_s$ ?


$s_1$     c b c a b a a

$s_2$     c b c a b c b

  
 $s_4$     c c c a b c a

$s_5$     c b c a a c a

$s_6$     c b c a a c a

  
 $s_8$     b b c a a c a

---

$s$     c b c a b c a

$k = 8$

$\ell = 7$

$t = 2$

$d_r = 2$

$d_s = 8$

# The “Outlier Variant”

## $(r, s)$ -CLOSEST STRING WITH OUTLIERS

*Instance:* Multi-set  $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^\ell$ ,  $\ell \in \mathbb{N}$ ,  
 $d_r, d_s, t \in \mathbb{N}$ .

*Question:* Is there an  $s \in \Sigma^\ell$ ,  
 $S' \subseteq S$  with  $|S'| = k - t$  with  
 $r_H(s, S') \leq d_r$  and  $s_H(s, S') \leq d_s$ ?

$s_1$     c b c a b a a

$s_2$     c b c a b c b

$s_4$     c c c a b c a

$s_5$     c b c a a c a

$s_6$     c b c a a c a

$s_8$     b b c a a c a

$s$     c b c a b c a

$$k = 8$$

$$\ell = 7$$

$$t = 2$$

$$d_r = 2$$

$$d_s = 8$$

$$r_H(s, S) = 2$$

$$s_H(s, S) = 7$$

## The “r- and s-Variants”

(r)-CLOSEST STRING:                    (r, s)-CLOSEST STRING **without**  $d_s$  bound  
(s)-CLOSEST STRING:                    (r, s)-CLOSEST STRING **without**  $d_r$  bound

Likewise for substring- and outlier-variants.

## The “r- and s-Variants”

(r)-CLOSEST STRING:	(r, s)-CLOSEST STRING without $d_s$ bound
(s)-CLOSEST STRING:	(r, s)-CLOSEST STRING without $d_r$ bound

Likewise for substring- and outlier-variants.

The (r)- and (s)-variants are intensively investigated:

our terminology	common in literature
(r)-CLOSEST STRING	CLOSEST STRING
(r)-CLOSEST SUBSTRING	CLOSEST SUBSTRING
(s)-CLOSEST SUBSTRING	CONSENSUS PATTERNS



## The “r- and s-Variants”

(r)-CLOSEST STRING:	(r, s)-CLOSEST STRING without $d_s$ bound
(s)-CLOSEST STRING:	(r, s)-CLOSEST STRING without $d_r$ bound

Likewise for substring- and outlier-variants.

The (r)- and (s)-variants are intensively investigated:

our terminology	common in literature
(r)-CLOSEST STRING	CLOSEST STRING
(r)-CLOSEST SUBSTRING	CLOSEST SUBSTRING
(s)-CLOSEST SUBSTRING	CONSENSUS PATTERNS

### Hardness

All these problems are **NP-hard**  
(except (s)-CLOSEST STRING, which is trivial).

## Parameters

$k$	number of input strings
$\ell$	length of input strings
$d_r$	radius bound
$d_s$	distance sum bound
$ \Sigma $	alphabet size
$m$	substring length (( $r, s$ )-CLOSEST SUBSTRING)
$t$	number of outliers (( $r, s$ )-CLOSEST STRING-WO)
$k - t$	number of inliers (( $r, s$ )-CLOSEST STRING-WO)

## Parameters

$k$	number of input strings
$\ell$	length of input strings
$d_r$	radius bound
$d_s$	distance sum bound
$ \Sigma $	alphabet size
$m$	substring length (( $r, s$ )-CLOSEST SUBSTRING)
$t$	number of outliers (( $r, s$ )-CLOSEST STRING-WO)
$k - t$	number of inliers (( $r, s$ )-CLOSEST STRING-WO)

Notation:

( $r, s$ )-CLOSEST STRING( $p_1, p_2, \dots$ ) means ( $r, s$ )-CLOSEST STRING parameterised by  $p_1, p_2, \dots$ .

## Parameters

$k$	number of input strings
$\ell$	length of input strings
$d_r$	radius bound
$d_s$	distance sum bound
$ \Sigma $	alphabet size
$m$	substring length (( $r, s$ )-CLOSEST SUBSTRING)
$t$	number of outliers (( $r, s$ )-CLOSEST STRING-WO)
$k - t$	number of inliers (( $r, s$ )-CLOSEST STRING-WO)

Notation:

( $r, s$ )-CLOSEST STRING( $p_1, p_2, \dots$ ) means ( $r, s$ )-CLOSEST STRING parameterised by  $p_1, p_2, \dots$ .

fixed-parameter tractable:  $\exists$  algorithm with running time  $f(k) \times p(|x|)$  for recursive  $f$  and polynomial  $p$  ( $x$  is input and  $k$  the parameter).

W[1]-hardness  $\Rightarrow$  *not* fixed parameter tractable.

# State of the Art

Previous literature:

(r)-CLOSEST STRING

(s)-CLOSEST STRING

(r)-CLOSEST SUBSTRING

(s)-CLOSEST SUBSTRING

((s)-CLOSEST SUBSTRING( $\ell, m$ ))

# State of the Art

Previous literature:

(r)-CLOSEST STRING

(s)-CLOSEST STRING

(r)-CLOSEST SUBSTRING

(s)-CLOSEST SUBSTRING

((s)-CLOSEST SUBSTRING( $\ell, m$ ))

(r, s)-CLOSEST STRING

(r, s)-CLOSEST SUBSTRING

# State of the Art

Previous literature:

(r)-CLOSEST STRING

(s)-CLOSEST STRING

(r)-CLOSEST SUBSTRING

(s)-CLOSEST SUBSTRING

((s)-CLOSEST SUBSTRING( $\ell, m$ ))

(r, s)-CLOSEST STRING

(r, s)-CLOSEST SUBSTRING

(r)-CLOSEST STRING-WO

(s)-CLOSEST STRING-WO

(r, s)-CLOSEST STRING-WO

# State of the Art

Previous literature:

(r)-CLOSEST STRING

(s)-CLOSEST STRING

(r)-CLOSEST SUBSTRING

(s)-CLOSEST SUBSTRING

((s)-CLOSEST SUBSTRING( $\ell, m$ ))

(r, s)-CLOSEST STRING

(r, s)-CLOSEST SUBSTRING

(r)-CLOSEST STRING-WO

(s)-CLOSEST STRING-WO

(r, s)-CLOSEST STRING-WO

Our Contribution:

(r, s)-CLOSEST STRING

(r, s)-CLOSEST SUBSTRING

(s)-CLOSEST SUBSTRING( $\ell, m$ )



# State of the Art

Previous literature:

(r)-CLOSEST STRING

(s)-CLOSEST STRING

(r)-CLOSEST SUBSTRING

(s)-CLOSEST SUBSTRING

((s)-CLOSEST SUBSTRING( $\ell, m$ ))

(r, s)-CLOSEST STRING

(r, s)-CLOSEST SUBSTRING

(r)-CLOSEST STRING-WO

(s)-CLOSEST STRING-WO

(r, s)-CLOSEST STRING-WO

Our Contribution:

(r, s)-CLOSEST STRING

(r, s)-CLOSEST SUBSTRING

(s)-CLOSEST SUBSTRING( $\ell, m$ )

(r)-CLOSEST STRING-WO

(s)-CLOSEST STRING-WO

(r, s)-CLOSEST STRING-WO

# Results for $(r, s)$ -CLOSEST STRING

$k$	$d_r$	$d_s$	$ \Sigma $	$\ell$	Result
<b>p</b>	–	–	–	–	FPT
–	<b>p</b>	–	–	–	FPT
–	–	<b>p</b>	–	–	FPT
–	–	–	2	–	NP-hard
–	–	–	–	<b>p</b>	FPT

# Results for $(r, s)$ -CLOSEST STRING

$k$	$d_r$	$d_s$	$ \Sigma $	$\ell$	Result
<b>p</b>	–	–	–	–	<b>FPT</b>
–	<b>p</b>	–	–	–	<b>FPT</b>
–	–	<b>p</b>	–	–	<b>FPT</b>
–	–	–	2	–	<b>NP-hard</b>
–	–	–	–	<b>p</b>	<b>FPT</b>

# Branching Algorithm

Fpt-branching algorithm for (r)-CLOSEST STRING( $d_r$ ):<sup>1</sup>

$$d_r = 2$$

$s_1$	c	b	c	a	b	a	a
$s_2$	c	b	c	a	b	c	b
$s_3$	a	b	c	c	c	c	a
$s_4$	c	c	c	a	b	c	a
$s_5$	c	b	c	a	a	c	a
$s_6$	c	b	c	a	a	c	a
$s_7$	a	b	b	a	b	a	a
$s_8$	b	b	c	a	a	c	a
<hr/>							
$s$	a	b	c	a	b	c	a

---

<sup>1</sup>Gramm, Niedermeier, Rossmanith, Algorithmica, 2003.

# Branching Algorithm

Fpt-branching algorithm for (r)-CLOSEST STRING( $d_r$ ):<sup>1</sup>

$$d_r = 2$$

$s_1$     c b c a b a a

$s_2$     c b c a b c b

$s_3$     a b c c c c a

$s_4$     c c c a b c a

$s_5$     c b c a a c a

$s_6$     c b c a a c a

$s_7$     a b b a b a a

$s_8$     b b c a a c a

---

$s$     a b c a b c a

$$s' = c b c a b a a$$

---

<sup>1</sup>Gramm, Niedermeier, Rossmanith, Algorithmica, 2003.

# Branching Algorithm

Fpt-branching algorithm for (r)-CLOSEST STRING( $d_r$ ):<sup>1</sup>

$$d_r = 2$$

$s_1$	c	b	c	a	b	a	a
$s_2$	c	b	c	a	b	c	b
$s_3$	a	b	c	c	c	c	a
$s_4$	c	c	c	a	b	c	a
$s_5$	c	b	c	a	a	c	a
$s_6$	c	b	c	a	a	c	a
$s_7$	a	b	b	a	b	a	a
$s_8$	b	b	c	a	a	c	a
<hr/>							
$s$	a	b	c	a	b	c	a

$$s' = \text{c b c a b a a}$$

---

<sup>1</sup>Gramm, Niedermeier, Rossmanith, Algorithmica, 2003.

## Extended Branching Algorithm

Goal: Extend branching algorithm to  $(\mathbf{r}, \mathbf{s})$ -CLOSEST STRING-WO( $d_r, t$ ).

## Extended Branching Algorithm

Goal: Extend branching algorithm to  $(r, s)$ -CLOSEST STRING-WO( $d_r, t$ ).

Choice of outliers can be added to the branching ( $t$  is a parameter).<sup>2</sup>

General branching similar:

branch by  $d_r + 1$  mismatches between candidate and input string.

---

<sup>2</sup>Boucher and Ma, BMC Bioinformatics, 2011.



## Extended Branching Algorithm

Goal: Extend branching algorithm to  $(r, s)$ -CLOSEST STRING-WO( $d_r, t$ ).

Choice of outliers can be added to the branching ( $t$  is a parameter).<sup>2</sup>

General branching similar:

branch by  $d_r + 1$  mismatches between candidate and input string.

Main problem: What is a good initial candidate string?

some input string  $\rightsquigarrow$  how do we satisfy  $d_s$  bound?

some string with low distance sum  $\rightsquigarrow$  how to bound depth of branching?

---

<sup>2</sup>Boucher and Ma, BMC Bioinformatics, 2011.

# Extended Branching Algorithm

Majority string  $s_m$ :

pick a most frequent symbol in each column.

Example:

$s_1$	d	b	a	b	b	b	b
$s_2$	d	a	a	b	c	c	d
$s_3$	d	a	a	b	c	c	d
$s_4$	a	a	c	c	c	c	d
$s_5$	a	a	c	b	c	c	d
$s_6$	a	c	a	b	d	b	d
$s_m$	d	a	a	b	c	c	d

# Extended Branching Algorithm

Majority string  $s_m$ :

pick a most frequent symbol in each column.

Example:

$s_1$	d	b	a	b	b	b	b
$s_2$	d	a	a	b	c	c	d
$s_3$	d	a	a	b	c	c	d
$s_4$	a	a	c	c	c	c	d
$s_5$	a	a	c	b	c	c	d
$s_6$	a	c	a	b	d	b	d
<hr/>							
$s_m$	d	a	a	b	c	c	d

## Lemma

$$r_H(s, S) \leq d_r \Rightarrow d_H(s_m, s) \leq 2d_r.$$

# Extended Branching Algorithm

Majority string  $s_m$ :

pick a most frequent symbol in each column.

Example:

$s_1$	d	b	a	b	b	b	b
$s_2$	d	a	a	b	c	c	d
$s_3$	d	a	a	b	c	c	d
$s_4$	a	a	c	c	c	c	d
$s_5$	a	a	c	b	c	c	d
$s_6$	a	c	a	b	d	b	d
<hr/>							
$s_m$	d	a	a	b	c	c	d

## Lemma

$$r_H(s, S) \leq d_r \Rightarrow d_H(s_m, s) \leq 2d_r.$$

Problem: Outliers!

# Extended Branching Algorithm

Majority string  $s_m$ :

pick a most frequent symbol in each column.

Example:

$t = 2$

$s_1$	d	b	a	b	b	b	b
$s_2$	d	a	a	b	c	c	d
$s_3$	d	a	a	b	c	c	d
$s_4$	a	a	c	c	c	c	d
$s_5$	a	a	c	b	c	c	d
$s_6$	a	c	a	b	d	b	d
<hr/>							
$s_m$	d	a	a	b	c	c	d

## Lemma

$$r_H(s, S) \leq d_r \Rightarrow d_H(s_m, s) \leq 2d_r.$$

Problem: Outliers!

# Extended Branching Algorithm

**Refined** majority string  $s_m^\diamond$ :

$\exists$  symb. with at least as many occ. as majority symbol minus  $t \Rightarrow$  use  $\diamond$ .

Example:

$t = 2$

$s_1$     d b a b b b b

$s_2$     d a a b c c d

$s_3$     d a a b c c d

$s_4$     a a c c c c d

$s_5$     a a c b c c d

$s_6$     a c a b d b d

---

$s_m^\diamond$      $\diamond$  a  $\diamond$  b c  $\diamond$  d

Lemma

$$r_H(s, S) \leq d_r \Rightarrow d_H(s_m, s) \leq 2d_r.$$

Problem: Outliers!

# Extended Branching Algorithm

**Refined** majority string  $s_m^\diamond$ :

$\exists$  symb. with at least as many occ. as majority symbol minus  $t \Rightarrow$  use  $\diamond$ .

Example:

$t = 2$

$s_1$	d	b	a	b	b	b	b
$s_2$	d	a	a	b	c	c	d
$s_3$	d	a	a	b	c	c	d
$s_4$	a	a	c	c	c	c	d
$s_5$	a	a	c	b	c	c	d
$s_6$	a	c	a	b	d	b	d
<hr/>							
$s_m^\diamond$	$\diamond$	a	$\diamond$	b	c	$\diamond$	d

Lemma

$$r_H(s, \mathcal{S}) \leq d_r \Rightarrow d_H(s_m, s) \leq 2d_r.$$

Problem: Outliers!

*Disputed columns.*

# Extended Branching Algorithm

**Refined** majority string  $s_m^\diamond$ :

$\exists$  symb. with at least as many occ. as majority symbol minus  $t \Rightarrow$  use  $\diamond$ .

Example:

$t = 2$

$s_1$	d	b	a	b	b	b	b
$s_2$	d	a	a	b	c	c	d
$s_3$	d	a	a	b	c	c	d
$s_4$	a	a	c	c	c	c	d
$s_5$	a	a	c	b	c	c	d
$s_6$	a	c	a	b	d	b	d
<hr/>							
$s_m^\diamond$	$\diamond$	a	$\diamond$	b	c	$\diamond$	d

**Lemma**

$$r_H(s, S) \leq d_r \Rightarrow d_H(s_m, s) \leq 2d_r.$$

Problem: Outliers!

*Disputed columns.*

**Lemma**

If the instance has a solution, then  $\#$ disp. columns  $\leq 4d_r$ .



# Extended Branching Algorithm

**Refined** majority string  $s_m^\diamond$ :

$\exists$  symb. with at least as many occ. as majority symbol minus  $t \Rightarrow$  use  $\diamond$ .

Example:

$t = 2$

$s_1$	d	b	a	b	b	b	b
$s_2$	d	a	a	b	c	c	d
$s_3$	d	a	a	b	c	c	d
$s_4$	a	a	c	c	c	c	d
$s_5$	a	a	c	b	c	c	d
$s_6$	a	c	a	b	d	b	d
<hr/>							
$s_m^\diamond$	$\diamond$	a	$\diamond$	b	c	$\diamond$	d

## Lemma

$$r_H(s, S) \leq d_r \Rightarrow d_H(s_m, s) \leq 2d_r.$$

Problem: Outliers!

*Disputed columns.*

## Lemma

If the instance has a solution, then #disp. columns  $\leq 4d_r$ .

Algorithm:

start with  $s_m^\diamond$

branch over  $d_r + 1$  mismatches (or declare outlier)

depth bound:  $6d_r + t$ .

# Results for $(r, s)$ -CLOSEST SUBSTRING

$\ell$	$k$	$m$	$d_r$	$d_s$	$ \Sigma $	Result
–	–	<b>p</b>	–	–	<b>p</b>	FPT
<b>p</b>	<b>p</b>	–	–	–	–	FPT
<b>p</b>	–	–	–	<b>p</b>	–	FPT
<b>p</b>	–	–	–	–	<b>p</b>	FPT
<b>p</b>	–	<b>p</b>	<b>p</b>	–	–	W[1]-hard
–	<b>p</b>	–	<b>p</b>	<b>p</b>	<b>p</b>	W[1]-hard
–	<b>p</b>	<b>p</b>	<b>p</b>	<b>p</b>	–	W[1]-hard

# Results for $(r, s)$ -CLOSEST SUBSTRING

$\ell$	$k$	$m$	$d_r$	$d_s$	$ \Sigma $	Result
–	–	<b>p</b>	–	–	<b>p</b>	FPT
<b>p</b>	<b>p</b>	–	–	–	–	FPT
<b>p</b>	–	–	–	<b>p</b>	–	FPT
<b>p</b>	–	–	–	–	<b>p</b>	FPT
<b>p</b>	–	<b>p</b>	<b>p</b>	–	–	W[1]-hard
–	<b>p</b>	–	<b>p</b>	<b>p</b>	<b>p</b>	W[1]-hard
–	<b>p</b>	<b>p</b>	<b>p</b>	<b>p</b>	–	W[1]-hard

## $(r, s)$ -CLOSEST SUBSTRING( $\ell, m$ )

### Theorem

$(s)$ -CLOSEST SUBSTRING( $\ell, m$ ) is  $\mathbf{W}[1]$ -hard.

## $(r, s)$ -CLOSEST SUBSTRING( $\ell, m$ )

### Theorem

$(s)$ -CLOSEST SUBSTRING( $\ell, m$ ) is  $\mathbf{W}[1]$ -hard.

Reduction from multi-coloured clique:

Let  $G = (V, E)$  be a graph with partition  $V = V_1 \cup \dots \cup V_{k_c}$  such that every  $V_i$ ,  $1 \leq i \leq k_c$ , is an independent set.

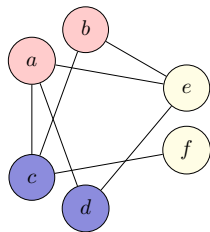
# $(r, s)$ -CLOSEST SUBSTRING( $\ell, m$ )

## Theorem

$(s)$ -CLOSEST SUBSTRING( $\ell, m$ ) is  $W[1]$ -hard.

Reduction from multi-coloured clique:

Let  $G = (V, E)$  be a graph with partition  $V = V_1 \cup \dots \cup V_{k_c}$  such that every  $V_i$ ,  $1 \leq i \leq k_c$ , is an independent set.



Repeat  $N = 36$  times

$\mathcal{V}_1$ :	\$	<b>a</b>	<b>c</b>	<b>e</b>	
$\mathcal{V}_2$ :	\$	<b>b</b>	<b>d</b>	<b>f</b>	
$\mathcal{E}_1$ :	\$	◇	<b>a</b>	<b>c</b>	◇
$\mathcal{E}_2$ :	\$	◇	<b>a</b>	<b>d</b>	◇
$\mathcal{E}_3$ :	\$	◇	<b>a</b>	◇	<b>e</b>
$\mathcal{E}_4$ :	\$	◇	<b>b</b>	<b>c</b>	◇
$\mathcal{E}_5$ :	\$	◇	<b>b</b>	◇	<b>e</b>
$\mathcal{E}_6$ :	\$	◇	◇	<b>c</b>	<b>f</b>
$\mathcal{E}_7$ :	\$	◇	◇	<b>d</b>	<b>e</b>
$s$ :	\$	<b>a</b>	<b>d</b>	<b>e</b>	

## Result for the Outlier Variants

### Theorem

(s)-CLOSEST STRING-WO( $d_s, \ell, k - t$ ) is W[1]-hard.

# Result for the Outlier Variants

## Theorem

(s)-CLOSEST STRING-WO( $d_s, \ell, k - t$ ) is  $W[1]$ -hard.

We know (r)-CLOSEST STRING-WO( $t = 0, |\Sigma| = 2$ ) is NP-hard, but ...

## Open Problem

(r)-CLOSEST STRING-WO( $|\Sigma|, k - t$ ),

(r)-CLOSEST STRING-WO( $|\Sigma|, d_r$ ),

(r)-CLOSEST STRING-WO( $|\Sigma|, d_r, k - t$ ).



## Result for the Outlier Variants

$(r, s)$ -CLOSEST STRING-WO( $\ell, |\Sigma|$ )  $\in$  FPT (trivial).

### Theorem

$(r, s)$ -CLOSEST STRING-WO( $\ell, |\Sigma|, d_r, d_s, (k - t)$ ) does not admit a polynomial kernel unless  $\text{coNP} \subseteq \text{NP/Poly}$ .

Thank you very much for your attention.