





The EXEMPLAR BREAKPOINT DISTANCE is not approximable for genomes with duplicated genes

Guillaume Blin¹ Guillaume Fertin² Florian Sikora¹ Stéphane Vialette¹

> ¹Université Paris-Est, LIGM - UMR CNRS 8049 - France {gblin,sikora,vialette}@univ-mlv.fr

²Université de Nantes, LINA - UMR CNRS 6241 - France guillaume.fertin@univ-nantes.fr

WALCOM February 2009

Florian Sikora

Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (1/23)

Outline

Motivations and notations

The EXEMPLAR BREAKPOINT DISTANCE (EBD) Problem

The ZERO EXEMPLAR BREAKPOINT DISTANCE (ZEBD) Problem

Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (2/23)

Outline

Motivations and notations

The EXEMPLAR BREAKPOINT DISTANCE (EBD) Problem

The ZERO EXEMPLAR BREAKPOINT DISTANCE (ZEBD) Problem

Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (3/23)

Motivations

Comparing two genomes (set of signed genes) species



Look for conserved set of genes in the same order



⊒

SQ (~

< ロ > < 国 > < 国 > < 国 > <

Look for a conserved set of genes

- Optimizing a given (dis)similarity measure
 - # breakpoints
 - # adjacencies
 - # conserved intervals
 - # common intervals
 - ▶ ...

Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (5/23)

Look for a conserved set of genes

Optimizing a given (dis)similarity measure

- # breakpoints
- # adjacencies
- # conserved intervals
- # common intervals
- ▶ ...

Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (5/23)

A non-conserved adjacency [Watterson et al. 1982]

Florian Sikora



Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (6/23)

<ロ> < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

SQ (V

A non-conserved adjacency [Watterson et al. 1982]





A non-conserved adjacency [Watterson et al. 1982]



Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (6/23)

A non-conserved adjacency [Watterson et al. 1982]



Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (6/23)

SQ (~

A non-conserved adjacency [Watterson et al. 1982]





A non-conserved adjacency [Watterson et al. 1982]



- Compute the number of breakpoints considering 2 genomes
- Done in polynomial time if there is no duplicated genes (only one occurrence of each gene family)
- Assumption which is not biologically valid

SQ (V

Exemplarization for dealing with duplicated genes

- Idea: perform an exemplarization which optimizes the number of breakpoints [Sankoff 1999]
- Exactly one occurrence of each gene family is kept (and matched)



The gene kept is assumed to be the ancestral (the other derived from it)

Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (7/23)

Outline

Motivations and notations

The EXEMPLAR BREAKPOINT DISTANCE (EBD) Problem

The ZERO EXEMPLAR BREAKPOINT DISTANCE (ZEBD) Problem

Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (8/23)

Definition of the EBD Problem

The EXEMPLAR BREAKPOINT DISTANCE Problem

Input: Two genomes G_1 , G_2 and an integer k**Question:** Is it possible to establish an exemplar matching of G_1 and G_2 such that the number of breakpoints between the resulting exemplar genomes is at most k ?

Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (9/23)

Previous results

G_1 +a -d +c -b -d -a +e +b -b

- occ(G₁) = 3 (maximum # occurences of a gene in the genome, here b)
- ► EBD(p,q): EBD problem with occ(G₁) = p and occ(G₂) = q

Previous results

G_1 +a -d +c -b -d -a +e +b -b

- occ(G₁) = 3 (maximum # occurences of a gene in the genome, here b)
- EBD(p,q): EBD problem with occ(G₁) = p and occ(G₂) = q
- ► *EBD*(1,2): **NP-Complete** [Bryant 2000]
- ► *EBD*(1,2): **APX-Hard** [Angibaud *et al.* 2008]
- ► *EBD*(3,3): **no approximation** algorithm [Chen *et al.* 2006]
- ▶ **Open** : Is *EBD*(2, 2) approximable ?

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ </p>

EBD Approximable ?

- OPT: the optimal solution of a problem Π
- A: a polynomial algorithm with a performance guarantee of

 α for each instance of Π
- $\mathcal{A} \leq \alpha \times OPT$

EBD Approximable ?

- OPT: the optimal solution of a problem Π
- A: a polynomial algorithm with a performance guarantee of

 α for each instance of Π
- $\mathcal{A} \leq \alpha \times OPT$
- *OPT* can be *zero* (if k = 0, no breakpoints allowed)
- If EBD with k = 0 is NP-Complete, no such algorithm A exists, Π is not approximable

Outline

Motivations and notations

The EXEMPLAR BREAKPOINT DISTANCE (EBD) Problem

The ZERO EXEMPLAR BREAKPOINT DISTANCE (ZEBD) Problem

Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (12/23)

Definition of the ZEBD Problem

The ZERO EXEMPLAR BREAKPOINT DISTANCE Problem

Input: Two genomes G_1 , G_2 **Question:** Is it possible to establish an exemplar matching of G_1 and G_2 such that the number of breakpoints between the resulting exemplar genomes is zero ?

Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (13/23)

Previous results

- ► ZEBD(3,3): NP-Complete [Chen et al. 2006]
- ZEBD(2, q): NP-Complete, q unbounded [Angibaud et al. 2008]
- ► **Open** : Is *ZEBD*(2, *q*) NP-Complete for bounded *q* ?

ZEBD is NP-Complete

- ► Our negative result: ZEBD(2,q) is **NP-Complete**
- Reduction from 3-SAT (proof in our article)

Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (15/23)

FPT for ZEBD

- Our positive result: ZEBD is Fixed Parameter Tractable (FPT)
- An FPT algorithm [Downey & Fellows 1999]: exact algorithm exponential only in its parameter (not in the input size)

- ► Parametrized by *m*, the number of genes families
- ► Using the **color-coding** technique by [Alon *et al.* 1995]



<ロ> < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Florian Sikora

- ► Parametrized by *m*, the number of genes families
- ► Using the **color-coding** technique by [Alon *et al.* 1995]



- ► Parametrized by *m*, the number of genes families
- Using the color-coding technique by [Alon et al. 1995]







- ► Parametrized by *m*, the number of genes families
- Using the color-coding technique by [Alon et al. 1995]



A vertex v of V for each pair of genes of the same family that carry the same sign

- ► Parametrized by *m*, the number of genes families
- ► Using the **color-coding** technique by [Alon *et al.* 1995]



A vertex v of V for each pair of genes of the same family that carry the same sign

- ► Parametrized by *m*, the number of genes families
- ► Using the **color-coding** technique by [Alon *et al.* 1995]



A vertex v of V for each pair of genes of the same family that carry the same sign

- ► Parametrized by *m*, the number of genes families
- Using the color-coding technique by [Alon et al. 1995]



A vertex v of V for each pair of genes of the same family that carry the same sign

- ► Parametrized by *m*, the number of genes families
- Using the color-coding technique by [Alon et al. 1995]



- A vertex v of V for each pair of genes of the same family that carry the same sign
- For all {(*i*, *j*), (*p*, *q*)} ∈ V², an edge from (*i*, *j*) to (*p*, *q*) if *i* < *p* and *j* < *q*

- ► Parametrized by *m*, the number of genes families
- Using the color-coding technique by [Alon et al. 1995]



- A vertex v of V for each pair of genes of the same family that carry the same sign
- For all {(*i*, *j*), (*p*, *q*)} ∈ V², an edge from (*i*, *j*) to (*p*, *q*) if *i* < *p* and *j* < *q*

- ► Parametrized by *m*, the number of genes families
- Using the color-coding technique by [Alon et al. 1995]



- A vertex v of V for each pair of genes of the same family that carry the same sign
- For all {(*i*, *j*), (*p*, *q*)} ∈ V², an edge from (*i*, *j*) to (*p*, *q*) if *i* < *p* and *j* < *q*

- Parametrized by *m*, the number of genes families
- Using the color-coding technique by [Alon et al. 1995]



- A vertex v of V for each pair of genes of the same family that carry the same sign
- For all {(*i*, *j*), (*p*, *q*)} ∈ V², an edge from (*i*, *j*) to (*p*, *q*) if *i* < *p* and *j* < *q*
- Looking for a colorful path: $\mathcal{O}(m2^m)$

The span: maximum distance between 2 occurrence of a gene family

G₁ +a -d +c -b -d -a +e +b -b



Florian Sikora

The span: maximum distance between 2 occurrence of a gene family



Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (18/23)

The span: maximum distance between 2 occurrence of a gene family

G₁ +a -d +c -b -d -a +e +b -b

Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (18/23)

The span: maximum distance between 2 occurrence of a gene family

G₁ +a -d +c -b -d -a +e +b -b

 Algorithm parametrized by the genome's span s (maximum span of a gene family in the genome)

•
$$span(G_1) = 5$$
 (for a family)

$$G_1 = \overbrace{+a-b+a} + c + d + c + d + e + f$$
$$G_2 = \underbrace{+a-b+c-b+d-f+e-f+e}$$

A vertex in the bag d_i for each common subsequence between each segment i of size s in G₁ and all G₂

Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (19/23)

$$G_1 = \overbrace{+a - b + a}^{\bullet} + c + d + c + d + e + f$$

 $G_2 = +a - b + c - b + d - f + e - f + e$
1 2 3 4 5 6 7 8 9

A vertex in the bag d_i for each common subsequence between each segment i of size s in G₁ and all G₂



Florian Sikora

Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (19/23)

$$G_1 = +a - b + a + c + d + c + d + e + f$$

 $G_2 = +a - b + c - b + d - f + e - f + e$
1 2 3 4 5 6 7 8 9

A vertex in the bag d_i for each common subsequence between each segment i of size s in G₁ and all G₂



Florian Sikora

Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (19/23)

$$G_1 = +a - b + a + c + d + c + d + e + f$$

 $G_2 = +a - b + c - b + d - f + e - f + e$
1 2 3 4 5 6 7 8 9

A vertex in the bag d_i for each common subsequence between each segment i of size s in G₁ and all G₂



Florian Sikora

Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (19/23)

$$G_1 = \overbrace{+a - b + a}^{+a - b + a} + c + d + c + d + e + f$$

 $G_2 = +a - b + c - b + d - f + e - f + e$
1 2 3 4 5 6 7 8 9

A vertex in the bag d_i for each common subsequence between each segment i of size s in G₁ and all G₂



 $v_1^1 = (+a, 1, 1)$ $v_1^2 = (-b, 2, 2)$ $v_1^3 = (-b, 4, 4)$ $v_1^4 = (+a - b, 1, 2)$

Florian Sikora

Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (19/23)

$$G_1 = +a - b + a + c + d + c + d + e + f$$

 $G_2 = +a - b + c - b + d - f + e - f + e$
1 2 3 4 5 6 7 8 9

A vertex in the bag d_i for each common subsequence between each segment i of size s in G₁ and all G₂



$$v_1^1 = (+a, 1, 1)$$

 $v_1^2 = (-b, 2, 2)$
 $v_1^3 = (-b, 4, 4)$
 $v_1^4 = (+a - b, 1, 2)$
 $v_1^5 = (+a - b, 1, 4)$

Florian Sikora Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (19/23)

$$G_1 = + a - b + a + c + d + c + d + e + f$$

 $G_2 = +a - b + c - b + d - f + e - f + e$
1 2 3 4 5 6 7 8 9

A vertex in the bag d_i for each common subsequence between each segment i of size s in G₁ and all G₂



$$v_1^1 = (+a, 1, 1)$$

 $v_1^2 = (-b, 2, 2)$
 $v_1^3 = (-b, 4, 4)$
 $v_1^4 = (+a - b, 1, 2)$
 $v_1^5 = (+a - b, 1, 4)$
 $v_2^1 = (+c, 3, 3)$

Florian Sikora

Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (19/23)

SQ (~

$$G_1 = + a - b + a + c + d + c + d + e + f$$

 $G_2 = +a - b + c - b + d - f + e - f + e$
1 2 3 4 5 6 7 8 9

- A vertex in the bag d_i for each common subsequence between each segment i of size s in G₁ and all G₂
- There is at most $|G_1| 2^s s$ vertices



$$v_1^1 = (+a, 1, 1)$$

 $v_1^2 = (-b, 2, 2)$
 $v_1^3 = (-b, 4, 4)$
 $v_1^4 = (+a - b, 1, 2)$
 $v_1^5 = (+a - b, 1, 4)$
 $v_2^1 = (+c, 3, 3)$

Florian Sikora

Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (19/23)

< □ > < □ > < 三 > < 三 > < 三 > < □ > < □ > <

$$G_1 = +a - b + a + c + d + c + d + e + f$$

 $G_2 = +a - b + c - b + d - f + e - f + e$
1 2 3 4 5 6 7 8 9

Add an edge if no gene of the same family in common (exemplarization) and not overlapping (no breakpoint)



$$v_1^1 = (+a, 1, 1)$$

 $v_1^2 = (-b, 2, 2)$
 $v_1^3 = (-b, 4, 4)$
 $v_1^4 = (+a - b, 1, 2)$
 $v_1^5 = (+a - b, 1, 4)$
 $v_2^1 = (+c, 3, 3)$

Florian Sikora

Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (20/23)

$$G_1 = +a - b + a + c + d + c + d + e + f$$

 $G_2 = +a - b + c - b + d - f + e - f + e$
1 2 3 4 5 6 7 8 9

 Add an edge if no gene of the same family in common (exemplarization) and not overlapping (no breakpoint)
 Ok



$$v_1^1 = (+a, 1, 1)$$

 $v_1^2 = (-b, 2, 2)$
 $v_1^3 = (-b, 4, 4)$
 $v_1^4 = (+a - b, 1, 2)$
 $v_1^5 = (+a - b, 1, 4)$
 $v_2^1 = (+c, 3, 3)$

Florian Sikora

Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (20/23)

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ </p>

$$G_1 = +a - b + a + c + d + c + d + e + f$$

 $G_2 = +a - b + c - b + d - f + e - f + e$
1 2 3 4 5 6 7 8 9

- Add an edge if no gene of the same family in common (exemplarization) and not overlapping (no breakpoint)
- Overlapping



$$v_1^1 = (+a, 1, 1)$$

 $v_1^2 = (-b, 2, 2)$
 $v_1^3 = (-b, 4, 4)$
 $v_1^4 = (+a - b, 1, 2)$
 $v_1^5 = (+a - b, 1, 4)$
 $v_2^1 = (+c, 3, 3)$

Florian Sikora

Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (20/23)

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ </p>

$$G_1 = +a -b + a + c + d + c + d + e + f$$

 $G_2 = +a -b + c - b + d - f + e - f + e$



Florian Sikora

Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (21/23)

Conclusion

- The EXEMPLAR BREAKPOINT DISTANCE problem cannot be approximated at all if both genomes contain duplicates
- The ZERO EXEMPLAR BREAKPOINT DISTANCE problem is NP-complete
- Two FPT algorithms for the ZERO EXEMPLAR BREAKPOINT DISTANCE problem
- Is there a constant ratio for EBD(1,q) ? (ZEBD(1,q) is polynomial, EBD(1,q) is APX-hard [Angibaud *et al.* 2008])

Questions on The EXEMPLAR BREAKPOINT DISTANCE is not approximable for genomes with duplicated genes ?

Guillaume Blin¹ Guillaume Fertin² Florian Sikora¹ Stéphane Vialette¹

> ¹Université Paris-Est, LIGM - UMR CNRS 8049 - France {gblin,sikora,vialette}@univ-mlv.fr

²Université de Nantes, LINA - UMR CNRS 6241 - France guillaume.fertin@univ-nantes.fr

WALCOM February 2009

Florian Sikora

Inapproximability of EXEMPLAR BREAKPOINT DISTANCE (23/23)