# Comparaisons de génomes avec gènes dupliqués : étude théorique et algorithmes

# Comparative genomics with duplicated genes: theoretical study and algorithms

Angibaud Sébastien

*sebastien.angibaud@univ-nantes.fr*

**L**aboratoire d'**I**nformatique de **N**antes **A**tlantique,
UMR CNRS 6241, UFR de Sciences et Techniques de Nantes

October 7th 2009



www.cnrs.fr

# Outline

1. Genomes comparison
   - Overview
   - Genomes representation
   - Measures between genomes

# Outline

# Outline

# Outline

# Outline

# Outline

# Genomes and genes

**Genome:**

- Composed of one or several *chromosomes*

# Genomes and genes

**Genome:**

- Composed of one or several *chromosomes*

# Genomes and genes

**Genome:**

- Composed of one or several *chromosomes*
- Sequence(s) of *DNA*
- Hereditary information

# Genomes and genes

**Genome:**

- Composed of one or several *chromosomes*
- Sequence(s) of *DNA*
- Hereditary information



**Gene:**

- Sequence of DNA
- Coding one or severals *proteins*
- Gene orientation

# Genomes and genes

**Genome:**

- Composed of one or several *chromosomes*
- Sequence(s) of *DNA*
- Hereditary information

**Gene:**

- Sequence of DNA
- Coding one or severals *proteins*
- Gene orientation

# Comparing genomes

**Why?**

# Comparing genomes

**Why?**

- Phylogenetic trees construction

# Comparing genomes

**Why?**

- Phylogenetic trees construction
- Identification of highly conserved sequences

# Comparing genomes

**Why?**

- Phylogenetic trees construction
- Identification of highly conserved sequences
- Help genome annotation

# Comparing genomes

**Why?**

- Phylogenetic trees construction
- Identification of highly conserved sequences
- Help genome annotation

**How?**

- Genome modeled as a sequence of genes

# Comparing two genomes : two different points of view

## Comparison based on the evolution process

- Infer an evolution process from one genome to another
- Several operations can be considered:
  - ▶ inversion
  - ▶ duplication
  - ▶ translocation
  - ▶ ...
- Find a most parsimonious rearrangement scenario

# Comparing two genomes : two different points of view

## Comparison based on the evolution process

- Infer an evolution process from one genome to another
- Several operations can be considered:
  - ▸ inversion
  - ▸ duplication
  - ▸ translocation
  - ▸ . . .
- Find a most parsimonious rearrangement scenario

## Comparison based on the structure of genomes

- Compare the structure (genes order) of the two genomes
- Compute a (dis)similarity measure between genomes
  - ▸ *number of breakpoints/adjacencies*
  - ▸ *number of common intervals*
  - ▸ *number of conserved intervals*
  - ▸ *Sum Adjacency Disruption*
  - ▸ *. . .*

# Comparing two genomes : two different points of view

## Comparison based on the evolution process

- Infer an evolution process from one genome to another
- Several operations can be considered:
  - ▶ inversion
  - ▶ duplication
  - ▶ translocation
  - ▶ . . .
- Find a most parsimonious rearrangement scenario

## Comparison based on the structure of genomes

- Compare the structure (genes order) of the two genomes
- Compute a (dis)similarity measure between genomes
  - ▶ *number of breakpoints/adjacencies*
  - ▶ *number of common intervals*
  - ▶ *number of conserved intervals*
  - ▶ *Sum Adjacency Disruption*
  - ▶ *. . .*

# Genomes representation

Representation and notations

1. Unichromosomal genome: sequence of *signed genes*

Example

1. $G_0 = +1 \ +2 \ -3 \ -7 \ +4 \ +5 \ +7 \ -8 \ +10 \ -9 \ +4 \ -6 \ -4$

# Genomes representation

### Representation and notations

1. Unichromosomal genome: sequence of *signed genes*
2. Alphabet $\Sigma \Leftrightarrow$ *gene families*

### Example

1. $G_0 = +1 \; +2 \; -3 \; \boxed{-7} \; +4 \; +5 \; \boxed{+7} \; -8 \; +10 \; -9 \; +4 \; -6 \; -4$
2. $\Sigma = \{1, 2, 3 \ldots 10\}$

# Genomes representation

## Representation and notations

1. Unichromosomal genome: sequence of *signed genes*
2. Alphabet $\Sigma \Leftrightarrow$ *gene families*
3. Let $G_0[k]$ be the $k^{th}$ gene (signed integer) of $G_0$

## Example

1. $G_0 = +1 \ +2 \ -3 \ -7 \ +4 \ +5 \ +7 \ -8 \ +10 \ -9 \ +4 \ -6 \ -4$
2. $\Sigma = \{1, 2, 3 \ldots 10\}$
3. $G_0[4] = -7$

# Genomes representation

### Representation and notations

1. Unichromosomal genome: sequence of *signed genes*
2. Alphabet $\Sigma \Leftrightarrow$ *gene families*
3. Let $G_0[k]$ be the $k^{th}$ gene (signed integer) of $G_0$
4. Let $occ(G_0)$ be the maximum number of genes in a gene family

### Example

1. $G_0 = +1 \; +2 \; -3 \; -7 \; +4 \; +5 \; +7 \; -8 \; +10 \; -9 \; +4 \; -6 \; -4$
2. $\Sigma = \{1, 2, 3 \dots 10\}$
3. $G_0[4] = -7$
4. $occ(G_0) = 3$

# Genomes representation

### Representation and notations

1. Unichromosomal genome: sequence of *signed genes*
2. Alphabet $\Sigma \Leftrightarrow$ *gene families*
3. Let $G_0[k]$ be the $k^{th}$ gene (signed integer) of $G_0$
4. Let $occ(G_0)$ be the maximum number of genes in a gene family
5. Let $\eta_{G_0}$ be the number of genes in $G_0$

### Example

1. $G_0 = +1 \ +2 \ -3 \ -7 \ +4 \ +5 \ +7 \ -8 \ +10 \ -9 \ +4 \ -6 \ -4$
2. $\Sigma = \{1, 2, 3 \ldots 10\}$
3. $G_0[4] = -7$
4. $occ(G_0) = 3$
5. $\eta_{G_0} = 13$

# Measures between two genomes

- **Input:** Two genomes $G_0$ and $G_1$ with the same gene contents and without duplicates
- **Output:** A (dis)-similarity measure between $G_0$ and $G_1$

- *number of breakpoints/adjacencies* [Watterson et al. 1982]
- *number of common intervals* [Uno and Yagiura, 2000]
- *number of conserved intervals* [Bergeron and Stoye, 2003]

# Breakpoint and adjacency

Definition: **adjacency** and **breakpoint** [Watterson et al. 1982]

There exists an adjacency between genes $G_0[p]$ and $G_0[p + 1]$ **iff**
$(G_0[p], G_0[p + 1])$ or $(-G_0[p + 1], -G_0[p])$ appears as a pair of
consecutive genes in $G_1$.

$$G_0 = +1 \ +2 \ +3 \ +4 \ +5$$
$$G_1 = +3 \ +4 \ -5 \ -2 \ -1$$

# Breakpoint and adjacency

Definition: **adjacency** and **breakpoint** [Watterson et al. 1982]

There exists an adjacency between genes $G_0[p]$ and $G_0[p + 1]$ **iff** $(G_0[p], G_0[p + 1])$ or $(-G_0[p + 1], -G_0[p])$ appears as a pair of consecutive genes in $G_1$.

$$\overbrace{+1 \ + 2}^{\textit{Adjacency}}$$
$$G_0 = +1 \ + 2 \ + 3 \ + 4 \ + 5$$
$$G_1 = -4 \ - 3 \ - 5 \ +1 \ + 2$$

# Breakpoint and adjacency

## Definition: **adjacency** and **breakpoint** [Watterson et al. 1982]

There exists an adjacency between genes $G_0[p]$ and $G_0[p+1]$ **iff** $(G_0[p], G_0[p+1])$ or $(-G_0[p+1], -G_0[p])$ appears as a pair of consecutive genes in $G_1$.

$$\overbrace{}^{\textit{Adjacency}} \overbrace{}^{\textit{Adjacency}}$$
$$G_0 = +1 \ + 2 \ +3 \ + 4 + 5$$
$$G_1 = -4 \ - 3 \ - 5 \ + 1 \ + 2$$

# Breakpoint and adjacency

## Definition: **adjacency** and **breakpoint** [Watterson et al. 1982]

There exists a breakpoint between genes $G_0[p]$ and $G_0[p+1]$ **iff** neither $(G_0[p], G_0[p+1])$ nor $(-G_0[p+1], -G_0[p])$ appears as a pair of consecutive genes in $G_1$.

$$G_0 = \overbrace{+1 \ +2}^{\textit{Adjacency}} \blacktriangledown \overbrace{+3 \ +4}^{\textit{Adjacency}} \blacktriangledown +5$$
$$G_1 = -4 \ -3 \ -5 \ +1 \ +2$$

# Breakpoint and adjacency

Definition: **adjacency** and **breakpoint** [Watterson et al. 1982]

There exists a breakpoint between genes $G_0[p]$ and $G_0[p+1]$ **iff** neither $(G_0[p], G_0[p+1])$ nor $(-G_0[p+1], -G_0[p])$ appears as a pair of consecutive genes in $G_1$.

$$G_0 = +0 \overbrace{+1 \ +2}^{\textit{Adjacency}} \overbrace{+3 \ +4}^{\textit{Adjacency}} +5 +6$$
$$G_1 = +0 \ -4 \ -3 \ -5 \ +1 \ +2 +6$$

# Breakpoint and adjacency

### Definition: **adjacency** and **breakpoint** [Watterson et al. 1982]

There exists a breakpoint between genes $G_0[p]$ and $G_0[p+1]$ **iff** neither $(G_0[p], G_0[p+1])$ nor $(-G_0[p+1], -G_0[p])$ appears as a pair of consecutive genes in $G_1$.

$$G_0 = +0 \overbrace{+1 + 2}^{\textit{Adjacency}} \overbrace{+3 + 4}^{\textit{Adjacency}} + 5 + 6$$
$$G_1 = +0 - 4 - 3 - 5 + 1 + 2 + 6$$

**Two measures:**

- *Number of adjacencies:* similarity
- *Number of breakpoints:* dissimilarity

# Common interval

Definition: **common interval** [Uno and Yagiura, 2000]

- A substring $s_0$ of $G_0$ is a *common interval* of $(G_0, G_1)$ if, in $G_1$, there is a substring $s_1$ such that $s_1$ is a permutation of $s_0$ (without taking signs into account)

$$G_0 = +1 + 2 + 3 + 4 + 5 \qquad G_1 = +2 - 4 + 3 + 5 + 1$$

# Common interval

Definition: **common interval** [Uno and Yagiura, 2000]

- A substring $s_0$ of $G_0$ is a *common interval* of $(G_0, G_1)$ if, in $G_1$, there is a substring $s_1$ such that $s_1$ is a permutation of $s_0$ (without taking signs into account)

$$G_0 = +1 + 2 \; \boxed{+3 +4 +5} \qquad G_1 = +2 \; \boxed{-4 +3 +5} + 1$$

$\Rightarrow s_0 = +3 + 4 + 5 \quad s_1 = -4 + 3 + 5$
Substring $s_0$ is a common interval of $(G_0, G_1)$.

# Common interval

Definition: **common interval** [Uno and Yagiura, 2000]

- A substring $s_0$ of $G_0$ is a *common interval* of $(G_0, G_1)$ if, in $G_1$, there is a substring $s_1$ such that $s_1$ is a permutation of $s_0$ (without taking signs into account)

$$G_0 = +1 + 2 \; \boxed{+3 +4 +5} \qquad G_1 = +2 \; \boxed{-4 +3 +5} + 1$$

$\Rightarrow s_0 = +3 + 4 + 5 \quad s_1 = -4 + 3 + 5$
Substring $s_0$ is a common interval of $(G_0, G_1)$.

- **Number of common intervals of $(G_0, G_1)$:**
  Similarity measure between two genomes

# Conserved interval

Definition: **conserved interval**
Proposed in [Bergeron and Stoye, 2003] for n permutations

- common interval
- same extremities  OR reversed extremities

$$G_0 = +0 + 1 + 2 + 3 + 4 + 5$$

$$G_1 = -4 - 3 - 5 + 0 - 1 + 2$$

# Conserved interval

Definition: **conserved interval**
Proposed in [Bergeron and Stoye, 2003] for n permutations

- common interval
- same extremities   OR   reversed extremities

$$G_0 = \text{+0 +1 +2} + 3 + 4 + 5$$

$$G_1 = -4 - 3 - 5 \ \text{+0 -1+2}$$

# Conserved interval

Definition: **conserved interval**
Proposed in [Bergeron and Stoye, 2003] for n permutations

- common interval
- same extremities  OR  reversed extremities

$$G_0 = +0 + 1 + 2 \;+3\;+4\; + 5$$

$$G_1 = \;-4\;-3\; - 5 + 0 - 1 + 2$$

# Conserved interval

Definition: **conserved interval**
Proposed in [Bergeron and Stoye, 2003] for n permutations

- common interval
- same extremities OR reversed extremities

$$G_0 = +0 + 1 + 2 + 3 + 4 + 5$$
$$G_1 = -4 - 3 - 5 + 0 - 1 + 2$$

- **Number of conserved intervals of ($G_0$, $G_1$):**
  Similarity measure between two genomes

# And with duplicates?

1. Choose a one-to-one correspondence $\mathcal{M}$ of genes (a matching)
2. Rename or remove genes according to $\mathcal{M}$
3. Compute the (dis)-similarity measure

# And with duplicates?

1. Choose a one-to-one correspondence $\mathcal{M}$ of genes (a matching)
2. Rename or remove genes according to $\mathcal{M}$
3. Compute the (dis)-similarity measure

*exemplar model (E)*
[Sankoff, 99]
one occurrence for each
gene family in $\mathcal{M}$

$G_0 = +0 +1 -2 -1 -3 +4$

$G_1 = +0 -1 +2 -1 -3 -1 +4$

# And with duplicates?

1. Choose a one-to-one correspondence $\mathcal{M}$ of genes (a matching)
2. Rename or remove genes according to $\mathcal{M}$
3. Compute the (dis)-similarity measure

*exemplar model ($E$)*
[Sankoff, 99]
one occurrence for each
gene family in $\mathcal{M}$

$G_0 = +0\ +1\ -2\ -1\ -3\ +4$

$G_1 = +0\ -1\ +2\ -1\ -3\ -1\ +4$

## And with duplicates?

1. Choose a one-to-one correspondence $\mathcal{M}$ of genes (a matching)
2. Rename or remove genes according to $\mathcal{M}$
3. Compute the (dis)-similarity measure

*exemplar model (**E**)*
[Sankoff, 99]
one occurrence for each
gene family in $\mathcal{M}$

$G_0^E = +0 +1 -2 -3 +4$

$G_1^E = +0 +2 -1 -3 +4$

# And with duplicates?

1. Choose a one-to-one correspondence $\mathcal{M}$ of genes (a matching)
2. Rename or remove genes according to $\mathcal{M}$
3. Compute the (dis)-similarity measure

*exemplar model ($\boldsymbol{E}$)*
[Sankoff, 99]
one occurrence for each
gene family in $\mathcal{M}$

$G_0^E = +0^{\blacktriangledown}+1-2^{\blacktriangledown}-3+4$

$G_1^E = +0+2-1-3+4$

$\text{Bkp}(G_0^E, G_1^E) = 2$

# And with duplicates?

- ① Choose a one-to-one correspondence $\mathcal{M}$ of genes (a matching)
- ② Rename or remove genes according to $\mathcal{M}$
- ③ Compute the (dis)-similarity measure

*exemplar model ($\textbf{E}$)*
[Sankoff, 99]
one occurrence for each
gene family in $\mathcal{M}$

$G_0^E = +0^{\blacktriangledown} +1 -2^{\blacktriangledown} -3 +4$

$G_1^E = +0 +2 -1 -3 +4$

$\mathrm{Bkp}(G_0^E, G_1^E) = 2$

*maximum matching model ($\textbf{M}$)*
[Tang & al, 03]
a maximum number of
occurrences in $\mathcal{M}$

$G_0 = +0 +1 -2 -1 -3 +4$

$G_1 = +0 -1 +2 -1 -3 -1 +4$

# And with duplicates?

1. Choose a one-to-one correspondence $\mathcal{M}$ of genes (a matching)
2. Rename or remove genes according to $\mathcal{M}$
3. Compute the (dis)-similarity measure

*exemplar model (**E**)*
[Sankoff, 99]
one occurrence for each
gene family in $\mathcal{M}$

$G_0^E = +0^\blacktriangledown +1 -2^\blacktriangledown -3 +4$

$G_1^E = +0 +2 -1 -3 +4$

$\mathrm{Bkp}(G_0^E, G_1^E) = 2$

*maximum matching model (**M**)*
[Tang & al, 03]
a maximum number of
occurrences in $\mathcal{M}$

$G_0 = +0 \; +1 \; -2 \; -1 \; -3 +4$

$G_1 = +0 \; -1 \; +2 \; -1 \; -3 \; -1 +4$

# And with duplicates?

1. Choose a one-to-one correspondence $\mathcal{M}$ of genes (a matching)
2. Rename or remove genes according to $\mathcal{M}$
3. Compute the (dis)-similarity measure

*exemplar model (E)*
[Sankoff, 99]
one occurrence for each
gene family in $\mathcal{M}$

$G_0^E = +0^\blacktriangledown +1 -2^\blacktriangledown -3 +4$

$G_1^E = +0 +2 -1 -3 +4$

$\text{Bkp}(G_0^E, G_1^E) = 2$

*maximum matching model (M)*
[Tang & al, 03]
a maximum number of
occurrences in $\mathcal{M}$

$G_0^M = +0 +1' -2 -1'' -3 +4$

$G_1^M = +0 -1' +2 -1'' -3 +4$

# And with duplicates?

1. Choose a one-to-one correspondence $\mathcal{M}$ of genes (a matching)
2. Rename or remove genes according to $\mathcal{M}$
3. Compute the (dis)-similarity measure

*exemplar model (**E**)*
[Sankoff, 99]
one occurrence for each
gene family in $\mathcal{M}$

$G_0^E = +0^\blacktriangledown +1 -2^\blacktriangledown -3 +4$

$G_1^E = +0 +2 -1 -3 +4$

$\mathrm{Bkp}(G_0^E, G_1^E) = 2$

*maximum matching model (**M**)*
[Tang & al, 03]
a maximum number of
occurrences in $\mathcal{M}$

$G_0^M = +0^\blacktriangledown +1'^\blacktriangledown -2^\blacktriangledown -1'' -3 +4$

$G_1^M = +0 -1' +2 -1'' -3 +4$

$\mathrm{Bkp}(G_0^M, G_1^M) = 3$

# And with duplicates?

1. Choose a one-to-one correspondence $\mathcal{M}$ of genes (a matching)
2. Rename or remove genes according to $\mathcal{M}$
3. Compute the (dis)-similarity measure

*exemplar model (E)*
[Sankoff, 99]
one occurrence for each
gene family in $\mathcal{M}$

*maximum matching model (M)*
[Tang & al, 03]
a maximum number of
occurrences in $\mathcal{M}$

*Intermediate model (I)*

For each gene family,
at least one gene is kept in $\mathcal{M}$

# Several possible matchings?

*maximum matching model (**M**)*
[Tang & al, 03]
a maximum number of occurrences in $\mathcal{M}$

$$G_0 = +0 \; \boxed{+1} \; - 2 \; \boxed{-1} \; - 3 \; + 4$$

$$G_1 = +0 \; \boxed{-1} \; + 2 \; \boxed{-1} \; - 3 \; -1 \; + 4$$

# Several possible matchings?

*maximum matching model (**M**)*
[Tang & al, 03]
a maximum number of occurrences in $\mathcal{M}$

$$G_0 = +0 \;\; +1 \;\; - 2 \;\; -1 \;\; - 3 \;\; + 4$$

$$G_1 = +0 \;\; -1 \;\; + 2 \;\; -1 \;\; - 3 \;\; -1 \;\; + 4$$

# Measure between genomes with duplicates

## Problem

- **Input:**
  - Two genomes $G_0$ and $G_1$
  - A model $X \in \{E, M, I\}$

- **Output:** Find a matching $\mathcal{M}$ which satisfies the model $X$, and which optimizes the measure between $G_0^X$ and $G_1^X$

# Measure between genomes with duplicates

## Problem

- **Input:**
  - Two genomes $G_0$ and $G_1$
  - A model $X \in \{E, M, I\}$
- **Output:** Find a matching $\mathcal{M}$ which satisfies the model $X$, and which optimizes the measure between $G_0^X$ and $G_1^X$

| measure | problem |
| --- | --- |
| common interval | $ICOM_X$ |
| conserved interval | $ICONS_X$ |
| breakpoint | $BD_X$ |
| adjacency | $ADJ_X$ |

# Measure between genomes with duplicates

## Problem

- **Input:**
  - Two genomes $G_0$ and $G_1$
  - A model $X \in \{E, M, I\}$
- **Question:** Are there $G_0^X$ and $G_1^X$ which satisfy the model $X$, and which imply no breakpoint ?

| measure | problem | |
|---|---|---|
| common interval | $ICOM_X$ | |
| conserved interval | $ICONS_X$ | |
| breakpoint | $BD_X$ | $ZBD_X$ |
| adjacency | $ADJ_X$ | |

# Outline

## What do we know?

|  | *exemplar model* | *maximum matching model* | *intermediate model* |
|---|---|---|---|
| $ICOM_X$ $ICONS_X$ | **NP**-Complete [Chauve et al.] (instance **(1, 2)**) | | |
| $BD_X$ | **NP**-Complete [Bryant] (instance **(1, 2)**) **NP**-Complete [Blin et al.] * | | |
| $ZBD_X$ | **NP**-Complete [Chen et al.] (instance **(3, 3)**) | ? | ? |

instance $(a, b) \Leftrightarrow occ(G_0) = a$ and $occ(G_1) = b$

* only one family contains several occurrences

# Definition

### $\alpha$-approximation and PTAS

- Let $P$ be an optimization problem
- Let $I$ be an instance of $P$
- A polynomial algorithm $A$ is an $\alpha$-approximation iff
  - If $P$ is a problem of minimization, then $A(I) \leqslant \alpha \cdot optimal(I)$
  - If $P$ is a problem of maximization, then $A(I) \geqslant \frac{1}{\alpha} \cdot optimal(I)$

# Definition

### $\alpha$-approximation and PTAS

- Let $P$ be an optimization problem
- Let $I$ be an instance of $P$
- A polynomial algorithm $A$ is an $\alpha$-approximation iff
  - If $P$ is a problem of minimization, then $A(I) \leqslant \alpha \cdot optimal(I)$
  - If $P$ is a problem of maximization, then $A(I) \geqslant \frac{1}{\alpha} \cdot optimal(I)$
- A polynomial algorithm $B$ is a *Polynomial-Time Approximation Scheme* (PTAS) iff $\forall \epsilon > 0$
  - If $P$ is a problem of minimization, then $B(I) \leqslant (1 + \epsilon) \cdot optimal(I)$
  - If $P$ is a problem of maximization, then $B(I) \geqslant \frac{1}{1+\epsilon} \cdot optimal(I)$

# Definition

## $\alpha$-**approximation** and **PTAS**

- Let $P$ be an optimization problem
- Let $I$ be an instance of $P$
- A polynomial algorithm $A$ is an $\alpha$-*approximation* iff
    - If $P$ is a problem of minimization, then $A(I) \leqslant \alpha \cdot optimal(I)$
    - If $P$ is a problem of maximization, then $A(I) \geqslant \frac{1}{\alpha} \cdot optimal(I)$
- A polynomial algorithm $B$ is a *Polynomial-Time Approximation Scheme* (PTAS) **iff** $\forall \epsilon > 0$
    - If $P$ is a problem of minimization, then $B(I) \leqslant (1 + \epsilon) \cdot optimal(I)$
    - If $P$ is a problem of maximization, then $B(I) \geqslant \frac{1}{1+\epsilon} \cdot optimal(I)$

## **APX**-Hard Class

- If a problem $P$ is **APX**-Hard then $P$ does not admit a **PTAS**

# New results

| | exemplar model | maximum matching model | intermediate model |
|---|---|---|---|
| $ICOM_X$ $ICONS_X$ | **NP**-Complete [Chauve et al.] (instance **(1, 2)**) **APX**-Hard (instance **(1, 2)**) * | | |
| $BD_X$ | **NP**-Complete [Bryant] (instance **(1, 2)**) **APX**-Hard (instance **(1, 2)**) * | **NP**-Complete [Blin et al.] | |
| $ZBD_X$ | **NP**-Complete [Chen et al.] (instance **(3, 3)**) (instance **(2, $k$)**) * [Blin et al.] (instance **(2, 2)**) | polynomial * | $ZBD_I \equiv$ $ZBD_E$ * |
| $ADJ_X$ | $ADJ_E \simeq BD_E$ * | $ADJ_M \simeq BD_M$ * | $ADJ_I \neq BD_I$ * |

# New results

|  | *exemplar model* | *maximum matching model* | *intermediate model* |
|---|---|---|---|
| $ICOM_X$ $ICONS_X$ | **NP**-Complete [Chauve et al.] (instance **(1, 2)**) **APX**-Hard (instance **(1, 2)**) | | |
| $BD_X$ | **NP**-Complete [Bryant] (instance **(1, 2)**) **NP**-Complete [Blin et al.] **APX**-Hard (instance **(1, 2)**) | | |
| $ZBD_X$ | **NP**-Complete [Chen et al.] (instance **(3, 3)**) (instance **(2, k)**) [Blin et al.] (instance **(2, 2)**) | polynomial | $ZBD_I \equiv ZBD_E$ |
| $ADJ_X$ | $ADJ_E \simeq BD_E$ | $ADJ_M \simeq BD_M$ | $ADJ_I \neq BD_I$ |

$A \simeq B$ : An optimal solution for $A$ is an optimal solution for $B$

$A \neq B$ : An optimal solution for $A$ *is not necessarily* an optimal solution for $B$

# New results

|  | *exemplar model* | *maximum matching model* | *intermediate model* |
|---|---|---|---|
| $ICOM_X$ $ICONS_X$ | **NP**-Complete [Chauve et al.] (instance **(1, 2)**) **APX**-Hard (instance **(1, 2)**) | | |
| $BD_X$ | **NP**-Complete [Bryant] (instance **(1, 2)**) | **NP**-Complete [Blin et al.] | |
| | **APX**-Hard (instance **(1, 2)**) | | |
| $ZBD_X$ | **NP**-Complete [Chen et al.] (instance **(3, 3)**) (instance **(2, k)**) [Blin et al.] (instance **(2, 2)**) | polynomial | $ZBD_I \equiv$ $ZBD_E$ |
| $ADJ_X$ | $ADJ_E \simeq BD_E$ | $ADJ_M \simeq BD_M$ | $ADJ_I \neq BD_I$ |

**A** $\simeq$ **B** : An optimal solution for **A** is an optimal solution for **B**

**A** $\neq$ **B** : An optimal solution for **A** *is not necessarily* an optimal solution for **B**

## New results

|  | *exemplar model* | *maximum matching model* | *intermediate model* |
|---|---|---|---|
| $ICOM_X$ $ICONS_X$ | **NP**-Complete [Chauve et al.] (instance **(1, 2)**) **APX**-Hard (instance **(1, 2)**) | | |
| $BD_X$ | **NP**-Complete [Bryant] (instance **(1, 2)**) | **NP**-Complete [Blin et al.] | |
| | **APX**-Hard (instance **(1, 2)**) | | |
| $ZBD_X$ | **NP**-Complete [Chen et al.] (instance **(3, 3)**) (instance **(2, $k$)**) [Blin et al.] (instance **(2, 2)**) | polynomial | $ZBD_I \equiv$ $ZBD_E$ |
| $ADJ_X$ | $ADJ_E \simeq BD_E$ | $ADJ_M \simeq BD_M$ | $ADJ_I \neq BD_I$ |

$\Rightarrow$ Bad news : $ICOM_X$, $ICONS_X$ and $BD_X$ do not admit a polynomial-time approximation scheme (PTAS)

# New results

| | exemplar model | maximum matching model | intermediate model |
|---|---|---|---|
| $ICOM_X$ $ICONS_X$ | **NP**-Complete [Chauve et al.] (instance **(1, 2)**) **APX**-Hard (instance **(1, 2)**) | | |
| $BD_X$ | **NP**-Complete [Bryant] (instance **(1, 2)**) | **NP**-Complete [Blin et al.] | |
| | **APX**-Hard  (instance **(1, 2)**) | | |
| $ZBD_X$ | **NP**-Complete [Chen et al.] (instance **(3, 3)**) (instance **(2, $k$)**) [Blin et al.] (instance **(2, 2)**) | polynomial | $ZBD_I \equiv$ $ZBD_E$ |
| $ADJ_X$ | $ADJ_E \simeq BD_E$ | $ADJ_M \simeq BD_M$ | $ADJ_I \neq BD_I$ |

$\Rightarrow$ Bad news : $BD_E$ and $BD_I$ do not admit any $\alpha$-approximation, unless **P** = **NP**

Angibaud Sébastien                     Phd Thesis - Defense                     October 7th 2009     19 / 49

# New results

|  | *exemplar model* | *maximum matching model* | *intermediate model* |
|---|---|---|---|
| $ICOM_X$ $ICONS_X$ | **NP**-Complete [Chauve et al.] (instance **(1, 2)**) | | |
| | **APX**-Hard (instance **(1, 2)**) | | |
| $BD_X$ | **NP**-Complete [Bryant] (instance **(1, 2)**) | | |
| | | **NP**-Complete [Blin et al.] | |
| | **APX**-Hard  (instance **(1, 2)**) | | |
| $ZBD_X$ | **NP**-Complete [Chen et al.] (instance **(3, 3)**) (instance **(2, k)**) [Blin et al.] (instance **(2, 2)**) | polynomial | $ZBD_I \equiv$ $ZBD_E$ |
| $ADJ_X$ | $ADJ_E \simeq BD_E$ | $ADJ_M \simeq BD_M$ | $ADJ_I \neq BD_I$ |

$\Rightarrow$ Good news : $BD_M$ **could** admit an $\alpha$-approximation

# Outline

# Exact algorithm

### Problem

- **Input:**
  - Two genomes $G_0$ and $G_1$
  - A model $X \in \{E, M, I\}$
- **Output:** Find a matching $\mathcal{M}$ which satisfies the model $X$, and which optimizes the measure between $G_0^X$ and $G_1^X$

**Idea:** transformation into a pseudo boolean linear problem

# Pseudo-boolean linear problem

### Definition

- **Variables**: domain $= \{0, 1\}$
- **Constraints**: inequalities between weighted sum of variables
- **Objective function**: weighted sum of variables

### Example

- **Variables**: $x \in \{0, 1\}, y \in \{0, 1\}, z \in \{0, 1\}$
- **Constraints**:
    - $x + 2 \cdot y \geqslant 2$
    - $z + y \leqslant 1$
- **Objective function**:
  maximize $x + 2 \cdot y - z$

# Pseudo-boolean linear problem

## Definition

- **Variables**: boolean
- **Constraints**: inequalities between weighted sum of variables
- **Objective function**: weighted sum of variables

## Example

- **Variables**: $x \in \{0, 1\}, y \in \{0, 1\}, z \in \{0, 1\}$
- **Constraints**:
  - $x + 2 \cdot y \geqslant 2$
  - $z + y \leqslant 1$
- **Objective function**:
  maximize $x + 2 \cdot y - z$

# Pseudo-boolean linear problem

## Definition

- **Variables**: boolean
- **Constraints**: inequalities between weighted sum of variables
- **Objective function**: weighted sum of variables

## Example

- **Variables**: $x \in \{0, 1\}, y \in \{0, 1\}, z \in \{0, 1\}$
- **Constraints**:
  - $x + 2 \cdot y \geqslant 2$
  - $z + y \leqslant 1$
- **Objective function**:
  maximize $x + 2 \cdot y - z$

# Pseudo-boolean linear problem

## Definition
- **Variables**: boolean
- **Constraints**: inequalities between weighted sum of variables
- **Objective function**: weighted sum of variables

## Example
- **Variables**: $x \in \{0, 1\}, y \in \{0, 1\}, z \in \{0, 1\}$
- **Constraints**:
  - $x + 2 \cdot y \geqslant 2$
  - $z + y \leqslant 1$
- **Objective function**:
  maximize $x + 2 \cdot y - z$

$\Rightarrow$ Powerful solvers for this type of problem

# Transformation for $ICOM_E$: variables

- Variables $x$ and $l$:

# Transformation for $ICOM_E$: variables

- Variables $x$ and $l$:



$x_b^a$ true $\Leftrightarrow$ gene $G_0[a]$ and $G_1[b]$ are matched

# Transformation for $ICOM_E$: variables

- Variables $x$ and $l$:



$l_{k,l,m,n}$ true $\Leftrightarrow$ $[k, l]$ in $G_0$ is a common interval of $(G_0, G_1)$, and $[m, n]$ in $G_1$ is a permutation of $[k, l]$

# Transformation for $ICOM_E$: constraints

**Exemplar model**:
for each genome, only one occurrence of each gene family

C1: $\forall f \in \mathcal{F}_{G_0} \cup \mathcal{F}_{G_1}, \quad \displaystyle\sum_{\substack{1 \leqslant a \leqslant \eta_{G_0} \\ G_0[a]=f}} \quad \sum_{\substack{1 \leqslant b \leqslant \eta_{G_1} \\ G_1[b]=f}} x_b^a = 1$

# Transformation for *ICOM$_E$*: constraints



Validity of variables $I_{k,l,m,n}$

$$I_{k,\ell,m,n} + x_2^3 \leqslant 1$$

# Transformation for $ICOM_E$

Objective function:

$$\text{Maximize} \sum_{k,l,m,n} I_{k,l,m,n}$$

# Transformation for $ICOM_E$

**Variables:**

$\mathcal{I} = \{ l_{k,l,m,n} : 1 \leqslant k \leqslant \ell \leqslant \eta_{G_0} \wedge 1 \leqslant m \leqslant n \leqslant \eta_{G_1} \}$

$\mathcal{X} = \{ x_b^a : 1 \leqslant a \leqslant \eta_{G_0} \wedge 1 \leqslant b \leqslant \eta_{G_1} \wedge G_0[a] = G_1[b] \}$

**Constraints:**

(C.01) $\forall f \in \mathcal{F}_{G_0} \cup \mathcal{F}_{G_1}, \quad \displaystyle\sum_{\substack{1 \leqslant a \leqslant \eta_{G_0} \\ G_0[a]=f}} \sum_{\substack{1 \leqslant b \leqslant \eta_{G_1} \\ G_1[b]=f}} x_b^a = 1$

(C.02) $\forall l_{k,l,m,n} \in \mathcal{I}, \; \forall k < p < \ell, \; \forall 1 \leqslant r < m, \quad G_0[p] = G_1[r], \quad l_{k,l,m,n} + x_r^p \leqslant 1$

(C.03) $\forall l_{k,l,m,n} \in \mathcal{I}, \; \forall k < p < \ell, \; \forall n < r \leqslant \eta_{G_1}, \; G_0[p] = G_1[r], \quad l_{k,l,m,n} + x_r^p \leqslant 1$

(C.04) $\forall l_{k,l,m,n} \in \mathcal{I}, \; \forall m < r < n, \; \forall 1 \leqslant p < k, \quad G_0[p] = G_1[r], \quad l_{k,l,m,n} + x_r^p \leqslant 1$

(C.05) $\forall l_{k,l,m,n} \in \mathcal{I}, \; \forall m < r < n, \; \forall \ell < p \leqslant \eta_{G_0}, \; G_0[p] = G_1[r], \quad l_{k,l,m,n} + x_r^p \leqslant 1$

(C.06) $\forall l_{k,l,m,n} \in \mathcal{I}, \quad 4\, l_{k,l,m,n} - \displaystyle\sum_{\substack{m \leqslant r \leqslant n \\ G_0[k]=G_1[r]}} x_r^k - \sum_{\substack{m \leqslant s \leqslant n \\ G_0[\ell]=G_1[s]}} x_s^\ell - \sum_{\substack{k \leqslant p \leqslant \ell \\ G_0[p]=G_1[m]}} x_m^p - \sum_{\substack{k \leqslant q \leqslant \ell \\ G_0[q]=G_1[n]}} x_n^q \leqslant 0$

**Objective function:**

Maximize $\displaystyle\sum_{k,l,m,n} l_{k,l,m,n}$

# Transformation for $ICOM_E$

**Variables:**

$\mathcal{I} = \{I_{k,l,m,n} : 1 \leqslant k \leqslant \ell \leqslant \eta_{G_0} \wedge 1 \leqslant m \leqslant n \leqslant \eta_{G_1}\}$

$\mathcal{X} = \{x_b^a : 1 \leqslant a \leqslant \eta_{G_0} \wedge 1 \leqslant b \leqslant \eta_{G_1} \wedge G_0[a] = G_1[b]\}$

**Constraints:**

(C.01) $\forall f \in \mathcal{F}_{G_0} \cup \mathcal{F}_{G_1}, \quad \sum\limits_{\substack{1 \leqslant a \leqslant \eta_{G_0} \\ G_0[a]=f}} \sum\limits_{\substack{1 \leqslant b \leqslant \eta_{G_1} \\ G_1[b]=f}} x_b^a = 1$

(C.02) $\forall I_{k,l,m,n} \in \mathcal{I}, \ \forall k < p < \ell, \ \forall 1 \leqslant r < m, \quad G_0[p] = G_1[r], \quad I_{k,l,m,n} + x_r^p \leqslant 1$

(C.03) $\forall I_{k,l,m,n} \in \mathcal{I}, \ \forall k < p < \ell, \ \forall n < r \leqslant \eta_{G_1}, \ G_0[p] = G_1[r], \quad I_{k,l,m,n} + x_r^p \leqslant 1$

(C.04) $\forall I_{k,l,m,n} \in \mathcal{I}, \ \forall m < r < n, \ \forall 1 \leqslant p < k, \quad G_0[p] = G_1[r], \quad I_{k,l,m,n} + x_r^p \leqslant 1$

(C.05) $\forall I_{k,l,m,n} \in \mathcal{I}, \ \forall m < r < n, \ \forall \ell < p \leqslant \eta_{G_0}, \ G_0[p] = G_1[r], \quad I_{k,l,m,n} + x_r^p \leqslant 1$

(C.06) $\forall I_{k,l,m,n} \in \mathcal{I}, \quad 4\,I_{k,l,m,n} - \sum\limits_{\substack{m \leqslant r \leqslant n \\ G_0[k]=G_1[r]}} x_r^k - \sum\limits_{\substack{m \leqslant s \leqslant n \\ G_0[\ell]=G_1[s]}} x_s^\ell - \sum\limits_{\substack{k \leqslant p \leqslant \ell \\ G_0[p]=G_1[m]}} x_m^p - \sum\limits_{\substack{k \leqslant q \leqslant \ell \\ G_0[q]=G_1[n]}} x_n^q \leqslant 0$

**Objective function:**

Maximize $\sum\limits_{k,l,m,n} I_{k,l,m,n}$

# Transformation for $ICOM_E$

**Variables:**

$\mathcal{I} = \{ I_{k,l,m,n} : 1 \leqslant k \leqslant \ell \leqslant \eta_{G_0} \wedge 1 \leqslant m \leqslant n \leqslant \eta_{G_1} \}$

$\mathcal{X} = \{ x_b^a : 1 \leqslant a \leqslant \eta_{G_0} \wedge 1 \leqslant b \leqslant \eta_{G_1} \wedge G_0[a] = G_1[b] \}$

**Constraints:**

(C.01) $\forall f \in \mathcal{F}_{G_0} \cup \mathcal{F}_{G_1}$, $\displaystyle \sum_{\substack{1 \leqslant a \leqslant \eta_{G_0} \\ G_0[a]=f}} \sum_{\substack{1 \leqslant b \leqslant \eta_{G_1} \\ G_1[b]=f}} x_b^a = 1$

(C.02) $\forall I_{k,l,m,n} \in \mathcal{I}, \ \forall k < p < \ell, \ \forall 1 \leqslant r < m, \quad G_0[p] = G_1[r], \quad I_{k,l,m,n} + x_r^p \leqslant 1$

(C.03) $\forall I_{k,l,m,n} \in \mathcal{I}, \ \forall k < p < \ell, \ \forall n < r \leqslant \eta_{G_1}, \ G_0[p] = G_1[r], \quad I_{k,l,m,n} + x_r^p \leqslant 1$

(C.04) $\forall I_{k,l,m,n} \in \mathcal{I}, \ \forall m < r < n, \ \forall 1 \leqslant p < k, \quad G_0[p] = G_1[r], \quad I_{k,l,m,n} + x_r^p \leqslant 1$

(C.05) $\forall I_{k,l,m,n} \in \mathcal{I}, \ \forall m < r < n, \ \forall \ell < p \leqslant \eta_{G_0}, \ G_0[p] = G_1[r], \quad I_{k,l,m,n} + x_r^p \leqslant 1$

(C.06) $\forall I_{k,l,m,n} \in \mathcal{I}, \quad 4\, I_{k,l,m,n} - \displaystyle\sum_{\substack{m \leqslant r \leqslant n \\ G_0[k]=G_1[r]}} x_r^k - \sum_{\substack{m \leqslant s \leqslant n \\ G_0[\ell]=G_1[s]}} x_s^\ell - \sum_{\substack{k \leqslant p \leqslant \ell \\ G_0[p]=G_1[m]}} x_m^p - \sum_{\substack{k \leqslant q \leqslant \ell \\ G_0[q]=G_1[n]}} x_n^q \leqslant 0$

**Objective function:**

Maximize $\displaystyle \sum_{k,l,m,n} I_{k,l,m,n}$

# Transformation for $ICOM_E$

**Variables:**

$\mathcal{I} = \{I_{k,l,m,n} : 1 \leqslant k \leqslant \ell \leqslant \eta_{G_0} \wedge 1 \leqslant m \leqslant n \leqslant \eta_{G_1}\}$

$\mathcal{X} = \{x_b^a : 1 \leqslant a \leqslant \eta_{G_0} \wedge 1 \leqslant b \leqslant \eta_{G_1} \wedge G_0[a] = G_1[b]\}$

**Constraints:**

(C.01) $\forall f \in \mathcal{F}_{G_0} \cup \mathcal{F}_{G_1},\quad \displaystyle\sum_{\substack{1 \leqslant a \leqslant \eta_{G_0} \\ G_0[a]=f}} \sum_{\substack{1 \leqslant b \leqslant \eta_{G_1} \\ G_1[b]=f}} x_b^a = 1$

(C.02) $\forall I_{k,l,m,n} \in \mathcal{I},\ \forall k < p < \ell,\ \forall 1 \leqslant r < m,\quad G_0[p] = G_1[r],\quad I_{k,l,m,n} + x_r^p \leqslant 1$

(C.03) $\forall I_{k,l,m,n} \in \mathcal{I},\ \forall k < p < \ell,\ \forall n < r \leqslant \eta_{G_1},\ G_0[p] = G_1[r],\quad I_{k,l,m,n} + x_r^p \leqslant 1$

(C.04) $\forall I_{k,l,m,n} \in \mathcal{I},\ \forall m < r < n,\ \forall 1 \leqslant p < k,\quad G_0[p] = G_1[r],\quad I_{k,l,m,n} + x_r^p \leqslant 1$

(C.05) $\forall I_{k,l,m,n} \in \mathcal{I},\ \forall m < r < n,\ \forall \ell < p \leqslant \eta_{G_0},\ G_0[p] = G_1[r],\quad I_{k,l,m,n} + x_r^p \leqslant 1$

(C.06) $\forall I_{k,l,m,n} \in \mathcal{I},\quad 4\,I_{k,l,m,n} - \displaystyle\sum_{\substack{m \leqslant r \leqslant n \\ G_0[k]=G_1[r]}} x_r^k - \sum_{\substack{m \leqslant s \leqslant n \\ G_0[\ell]=G_1[s]}} x_s^\ell - \sum_{\substack{k \leqslant p \leqslant \ell \\ G_0[p]=G_1[m]}} x_m^p - \sum_{\substack{k \leqslant q \leqslant \ell \\ G_0[q]=G_1[n]}} x_n^q \leqslant 0$

**Objective function:**

Maximize $\displaystyle\sum_{k,l,m,n} I_{k,l,m,n}$

# Pseudo boolean transformation

### Other problems ?

- **other models:** modify constraints C1
- **conserved intervals:** restriction on variables $I_{k,\ell,m,n}$
- **breakpoint and adjacency:** new variables and constraints

- $ICOM_X$ **and** $ICONS_X$
  S. Angibaud, G. Fertin, I. Rusu et S. Vialette.
  A pseudo-boolean general framework for computing rearrangement distances between genomes with duplicates
  *Journal of Computational Biology*, Vol. 14(4), pages 379-393. 2007

- $BD_X$ **and** $ADJ_X$
  S. Angibaud, G. Fertin, I. Rusu, A. Thévenin et S. Vialette.
  Efficient Tools for Computing the Number of Breakpoints and the Number of Adjacencies between two Genomes with Duplicate Genes
  *Journal of Computational Biology*, Vol. 15(8), pages 1093-1115. 2008

# Experimental results

### Dataset

- Twelve genomes of $\gamma$-*Proteobacteria* [Lerat et al. 2003]

| Name | Genbank identifier | size |
|------|:------------------:|:----:|
| *Buchnera aphidicola APS* | NC_002528 | 564 |
| *Escherichia coli K12* | NC_000913 | 4183 |
| *Haemophilus influenzae Rd* | NC_000907 | 1709 |
| *Pseudomonas aeruginosa PA01* | NC_002516 | 5540 |
| *Pasteurella multocida Pm70* | NC_002663 | 2015 |
| *Salmonella typhimurium LT2* | NC_003197 | 4203 |
| *Wigglesworthia glossinidia brevipalpis* | NC_004344 | 653 |
| *Xanthomonas axonopodis pv. citri 306* | NC_003919 | 4192 |
| *Xanthomonas campestris* | NC_0 03902 | 4029 |
| *Xylella fastidiosa 9a5c* | NC_002488 | 2680 |
| *Yersinia pestis CO_92* | NC_003143 | 3599 |
| *Yersinia pestis KIM5 P12* | NC_004088 | 3879 |
| | average: | 3104 |

# Experimental results

### Dataset

- Twelve genomes of $\gamma$-*Proteobacteria* [Lerat et al. 2003]
- 66 possible pairs of genomes

**Number of results:**

|          | model    |                  |              |
|----------|----------|------------------|--------------|
|          | Exemplar | maximum matching | intermediate |
| $ADJ_X$  | 61/66    | 66/66            | 63/66        |
| $ICOM_X$ | 21/66    | 40/66            | 21/66        |

# Experimental results

## Dataset

- Twelve genomes of $\gamma$-*Proteobacteria* [Lerat et al. 2003]
- 66 possible pairs of genomes

**Number of results:**

|        |          | model            |              |
|--------|----------|------------------|--------------|
|        | Exemplar | maximum matching | intermediate |
| $ADJ_X$  | 61/66    | 66/66            | 63/66        |
| $ICOM_X$ | 21/66    | 40/66            | 21/66        |

$\Rightarrow$ Efficient approach for $ADJ_X$

# Experimental results

## Dataset

- Twelve genomes of $\gamma$-*Proteobacteria* [Lerat et al. 2003]
- 66 possible pairs of genomes

**Number of results:**

|  | model | | |
|---|---|---|---|
|  | Exemplar | maximum matching | intermediate |
| $ADJ_X$ | 61/66 | 66/66 | 63/66 |
| $ICOM_X$ | 21/66 | 40/66 | 21/66 |

- $\Rightarrow$ Efficient approach for $ADJ_X$
- $\Rightarrow$ Limit is attained for $ICOM_X$
  - $\Rightarrow$ Heuristics

# Outline

# IILCS$_M$ heuristic

- Based on ILCS$_M$ heuristic [Tichy, 82]
- **Idea:** Match genes of a Longest Common Substring (LCS)

# IILCS*M* heuristic

- Based on ILCS*M* heuristic [Tichy, 82]
- **Idea:** Match genes of a Longest Common Substring (LCS)

### IILCS*M* heuristic

1. Compute the Longest Common Substring **S**

### Example

$$+1 +2 +3 \ \ +4 \ +5 \ +6 \ +7$$

$$+6 \ -7 \ +4 \ +5 \ +1 \ +6 \ \text{-3 -2 -1}$$

# IILCS*M* heuristic

- Based on ILCS*M* heuristic [Tichy, 82]
- **Idea:** Match genes of a Longest Common Substring (LCS)

### IILCS*M* heuristic

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly

### Example

$$+1\ +2\ +3\ \ +4\ +5\ +6\ +7$$

$$+6\ -7\ +4\ +5\ +1\ +6\ \text{-3 -2 -1}$$

# IILCS*M* heuristic

- Based on ILCS*M* heuristic [Tichy, 82]
- **Idea:** Match genes of a Longest Common Substring (LCS)

## IILCS*M* heuristic

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Remove genes that cannot be matched

## Example

$$+1 \; +2 \; +3 \; + 4 \; + 5 \; + 6 \; + 7$$

$$+6 \; - 7 \; + 4 \; + 5 \; +1 \; + 6 \; -3 \; -2 \; -1$$

# IILCS$_M$ heuristic

- Based on ILCS$_M$ heuristic [Tichy, 82]
- **Idea:** Match genes of a Longest Common Substring (LCS)

## IILCS$_M$ heuristic

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Remove genes that cannot be matched

## Example

$$+1 \ +2 \ +3 \quad + \ 4 \ + \ 5 \ + \ 6 \ + \ 7$$

$$+6 \ - \ 7 \ + \ 4 \ + \ 5 \ + \ 6 \ \text{-3 -2 -1}$$

# IILCS$_M$ heuristic

- Based on ILCS$_M$ heuristic [Tichy, 82]
- **Idea:** Match genes of a Longest Common Substring (LCS)

## IILCS$_M$ heuristic

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Remove genes that cannot be matched
4. Iterate the process until saturation

## Example

$$+1 \; +2 \; +3 \quad + 4 \; + 5 \; + 6 \; + 7$$

$$+6 \; - 7 \; + 4 \; + 5 \; + 6 \quad -3 \; -2 \; -1$$

# IILCS*M* heuristic

- Based on ILCS*M* heuristic [Tichy, 82]
- **Idea:** Match genes of a Longest Common Substring (LCS)

## IILCS*M* heuristic

1. Compute the Longest Common Substring *S*
2. Match all the genes of *S* accordingly
3. Remove genes that cannot be matched
4. Iterate the process until saturation

## Example

+1 +2 +3  +4 +5 +6  + **7**

+**6** − **7** +4 +5 +6 -3 -2 -1

# IILCS*M* heuristic

- Based on ILCS*M* heuristic [Tichy, 82]
- **Idea:** Match genes of a Longest Common Substring (LCS)

## IILCS*M* heuristic

1. Compute the Longest Common Substring *S*
2. Match all the genes of *S* accordingly
3. Remove genes that cannot be matched
4. Iterate the process until saturation

## Example

$$+1 \ +2 \ +3 \quad +4 \ +5 \ +6 \quad + \ 7$$

$$+6 \quad - \ 7 \quad +4 \ +5 \ +6 \quad -3 \ -2 \ -1$$

# IILCS$_M$ heuristic

- Based on ILCS$_M$ heuristic [Tichy, 82]
- **Idea:** Match genes of a Longest Common Substring (LCS)

## IILCS$_M$ heuristic

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Remove genes that cannot be matched
4. Iterate the process until saturation

## Example

$$+1\ +2\ +3\quad +4\ +5\ +6\quad +\ 7$$
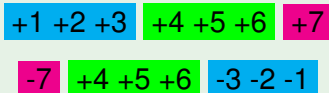
$$-7\quad +4\ +5\ +6\quad -3\ -2\ -1$$

# IILCS$_M$ heuristic

- Based on ILCS$_M$ heuristic [Tichy, 82]
- **Idea:** Match genes of a Longest Common Substring (LCS)

## IILCS$_M$ heuristic

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Remove genes that cannot be matched
4. Iterate the process until saturation

## Example

+1 +2 +3  +4 +5 +6  +7

-7  +4 +5 +6  -3 -2 -1

# IILCS*M* heuristic

- Based on ILCS*M* heuristic [Tichy, 82]
- **Idea:** Match genes of a Longest Common Substring (LCS)

## IILCS*M* heuristic

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Remove genes that cannot be matched
4. Iterate the process until saturation
5. Compute the measure

## Example

+1 +2 +3  +4 +5 +6  +7

-7  +4 +5 +6  -3 -2 -1

# Hybrid method

## Algorithm HYB$_X$(k)

- **Idea:** Associate exact method and IILCS$_X$ heuristic
- **Parameter $k$:** Bound on **LCS** size

1. Compute an **LCS $S$** of $(G_0, G_1)$
2. **If** $|S| \geqslant k$
   **Then**
   > Match all the genes of $S$
   > Remove genes that cannot be matched
   > Return to ①

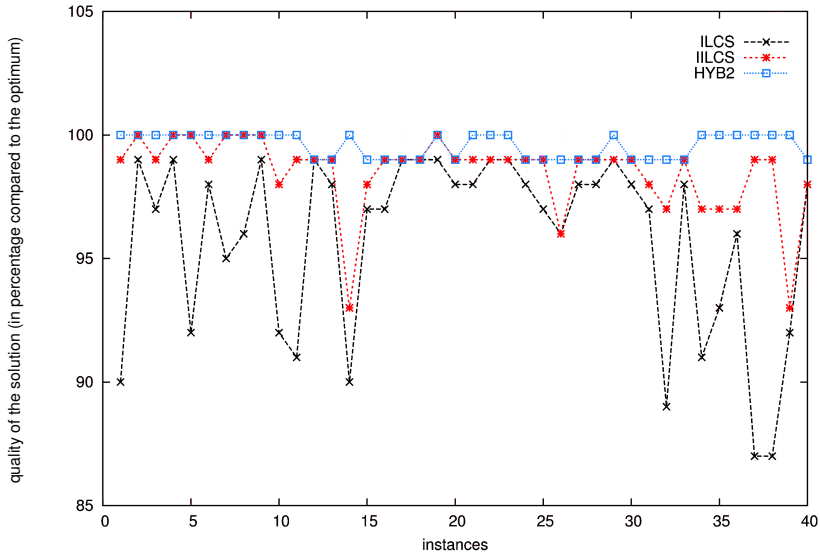   **Else** Apply the exact method: transformation into a pseudo-boolean linear problem
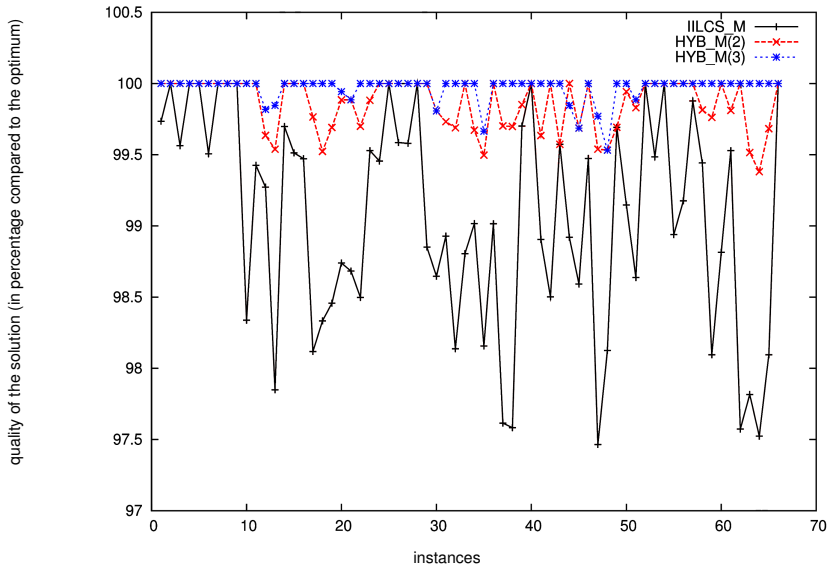
# Experimental results

## Dataset

- Twelve genomes of $\gamma$-*Proteobacteria* [Lerat et al. 2003]
- 66 possible pairs of genomes

| EXACT | model | | |
|---|---|---|---|
| | Exemplar | maximum matching | intermediate |
| $ADJ_X$ | 61/66 | 66/66 | 63/66 |
| $ICOM_X$ | 21/66 | 40/66 | 21/66 |

# Experimental results: $ICOM_M$

# Experimental results: $ADJ_M$

# Outline

# Goal

## Problem

- **Input:** two circular genomes $G_1$ and $G_2$
- **Output:** List of common intervals between $G_1$ and $G_2$

## Goal

- Compute common intervals
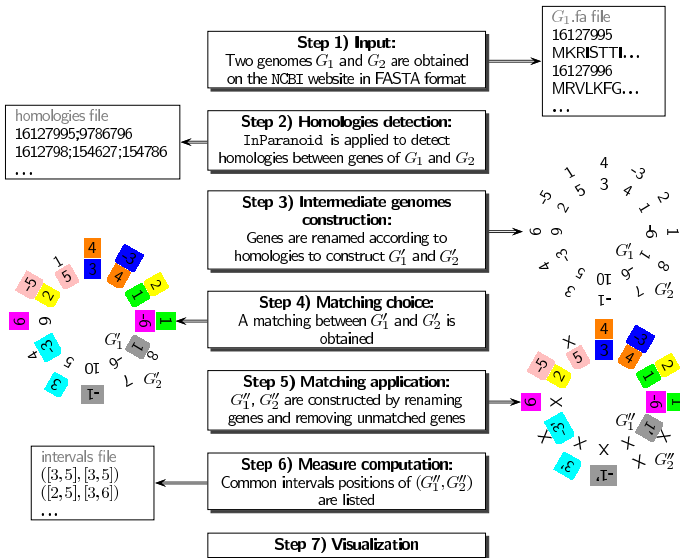- Provide a tool to visualize and analyze results

S. Angibaud, D. Éveillard, G. Fertin et I. Rusu
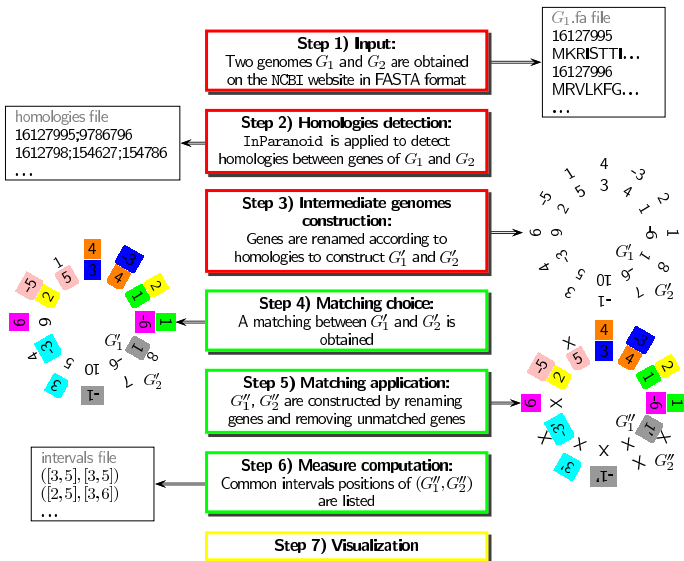Comparing Bacterial Genomes by Searching Their Common Intervals
*In Proc. 1st International Conference on Bioinformatics and Computational Biology*
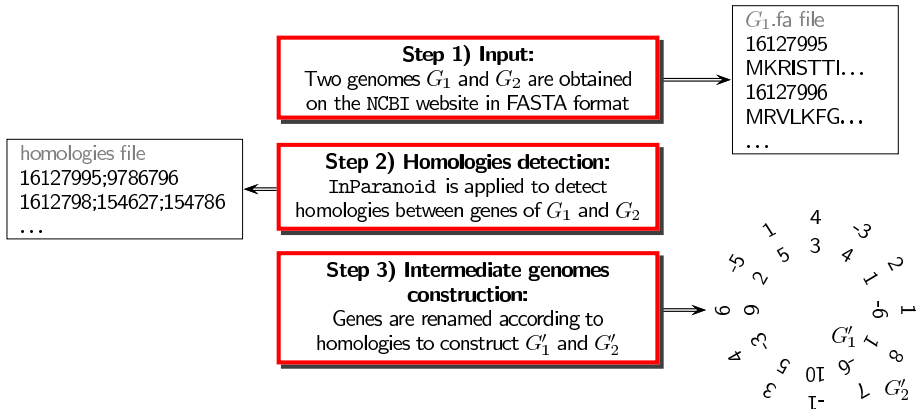LNBI Vol. 5462, pages 102-113. 2009

# Protocol

# Protocol

# Homologies computation



**Step 1) Input:**
Two genomes $G_1$ and $G_2$ are obtained
on the NCBI website in FASTA format

$G_1$.fa file
16127995
MKRISTTI...
16127996
MRVLKFG...
...

homologies file
16127995;9786796
1612798;154627;154786
...

**Step 2) Homologies detection:**
InParanoid is applied to detect
homologies between genes of $G_1$ and $G_2$

**Step 3) Intermediate genomes
construction:**
Genes are renamed according to
homologies to construct $G_1'$ and $G_2'$

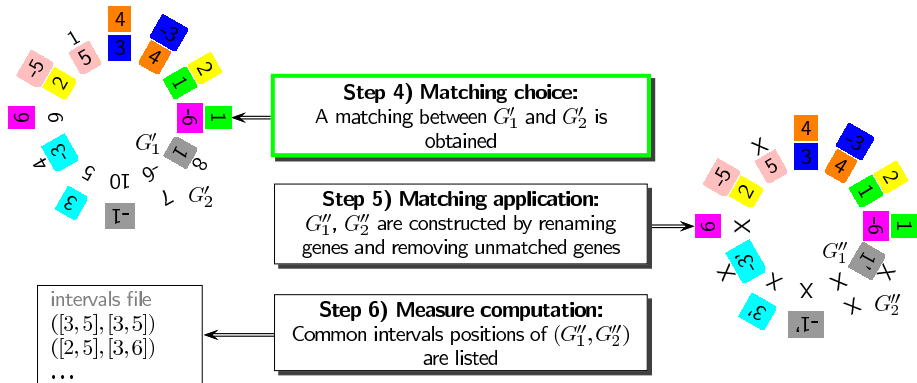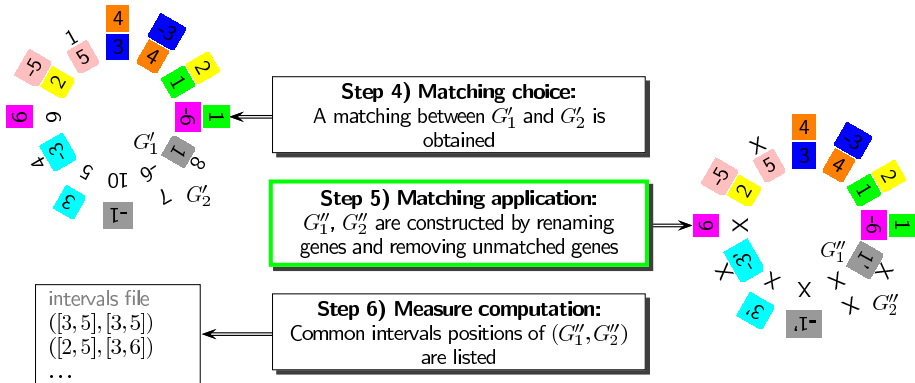## Inparanoid [Storm et al. 2001]

- Proposed in 2001 by Storm, Remm and Sonnhammer
- Compute clusters of homologous genes

# Step 4: choose a matching



- Exact method: Pseudo boolean transformation
- IILCS$_x$ heuristic
- Hybrid method

# Step 5: Matching application

# Step 6: common intervals computation



**Step 4) Matching choice:**
A matching between $G_1'$ and $G_2'$ is obtained

**Step 5) Matching application:**
$G_1''$, $G_2''$ are constructed by renaming genes and removing unmatched genes

**Step 6) Measure computation:**
Common intervals positions of $(G_1'', G_2'')$ are listed

intervals file
$([3, 5], [3, 5])$
$([2, 5], [3, 6])$
...

# Seven steps



**Step 1) Input:**
Two genomes $G_1$ and $G_2$ are obtained on the NCBI website in FASTA format

**Step 2) Homologies detection:**
InParanoid is applied to detect homologies between genes of $G_1$ and $G_2$

**Step 3) Intermediate genomes construction:**
Genes are renamed according to homologies to construct $G'_1$ and $G'_2$

**Step 4) Matching choice:**
A matching between $G'_1$ and $G'_2$ is obtained

**Step 5) Matching application:**
$G''_1$, $G''_2$ are constructed by renaming genes and removing unmatched genes

**Step 6) Measure computation:**
Common intervals positions of $(G''_1, G''_2)$ are listed

**Step 7) Visualization**

$G_1$.fa file
16127995
MKRISTTI...
16127996
MRVLKFG...
...

homologies file
16127995;9786796
1612798;154627;154786
...

intervals file
([3,5],[3,5])
([2,5],[3,6])
...

# Outline

# Contributions

- Better knowledge of problems
    - **APX**-Hardness of $BD_X$, $ICOM_X$ and $ICONS_X$
    - **NP**-Completeness of $ZBD_E$ and $ZBD_I$
    - Polynomiality of $ZBD_M$

# Contributions

- Better knowledge of problems
    - **APX**-Hardness of $BD_X$, $ICOM_X$ and $ICONS_X$
    - **NP**-Completeness of $ZBD_E$ and $ZBD_I$
    - Polynomiality of $ZBD_M$

- Three new algorithms
    - An exact approach based on a transformation into a pseudo-boolean problem
        - Efficient approach for $BD_X$ and $ADJ_X$
        - Limited for $ICOM_X$

# Contributions

- Better knowledge of problems
  - **APX**-Hardness of $BD_X$, $ICOM_X$ and $ICONS_X$
  - **NP**-Completeness of $ZBD_E$ and $ZBD_I$
  - Polynomiality of $ZBD_M$

- Three new algorithms
  - An exact approach based on a transformation into a pseudo-boolean problem
    - Efficient approach for $BD_X$ and $ADJ_X$
    - Limited for $ICOM_X$
  - IILCS$_X$ heuristic and Hybrid method
    - Promising results on a real dataset for each problem

# Perspectives

- Work on MATCH&WATCH
  - First experimentation on six chromosomes of $\gamma$-Proteobacteria
  - Analyze in details the common intervals obtained
  - Add functionalities according to biologists

# Perspectives

- Work on MATCH&WATCH
  - First experimentation on six chromosomes of $\gamma$-Proteobacteria
  - Analyze in details the common intervals obtained
  - Add functionalities according to biologists

- Multi-chromosomal genome comparison

- Multiple genome comparison

# Perspectives

- Work on MATCH&WATCH
  - First experimentation on six chromosomes of $\gamma$-Proteobacteria
  - Analyze in details the common intervals obtained
  - Add functionalities according to biologists

- Multi-chromosomal genome comparison

- Multiple genome comparison

- New algorithms
  - $\alpha$-approximation for $BD_E$ and $BD_I$ when $occ(G_0) = 1$?
  - $\alpha$-approximation or PTAS for $ICOM_X$ on balanced genomes?

# Perspectives

- Work on MATCH&WATCH
  - First experimentation on six chromosomes of $\gamma$-Proteobacteria
  - Analyze in details the common intervals obtained
  - Add functionalities according to biologists

- Multi-chromosomal genome comparison

- Multiple genome comparison

- New algorithms
  - $\alpha$-approximation for $BD_E$ and $BD_I$ when $occ(G_0) = 1$?
  - $\alpha$-approximation or PTAS for $ICOM_X$ on balanced genomes?

- Partially ordered genomes

# Acknowledgement

- *Directors*
  - ▶ Irena Rusu

  - ▶ Guillaume Fertin

- *Co-authors*
  - ▶ Damien Éveillard (LINA, Université de Nantes)

  - ▶ Annelyse Thévenin (LRI, Université Paris-Sud)

  - ▶ Stéphane Vialette (IGM, Université Paris-Est Marne-la-Vallée)

**Pictures**
- http://www.mun.ca/biology/scarr/FISH_chromosomes_300dpi.jpg
- http://agaudi.files.wordpress.com/2008/09/dna_overview_es.png
- http://joachimj.club.fr/imagesmada2004bis/PlanchePhylogeniedesprimates.jpg
- http://http://fr.wikipedia.org/wiki/Gene
- http://www.g-language.org/g3/

# Acknowledgement

- *Directors*
  - ▶ Irena Rusu

  - ▶ Guillaume Fertin

**Thank you**

**Merci**

- *Co-authors*
  - ▶ Damien Éveillard (LINA, Université de Nantes)

  - ▶ Annelyse Thévenin (LRI, Université Paris-Sud)

  - ▶ Stéphane Vialette (IGM, Université Paris-Est Marne-la-Vallée)

**Pictures**

  - http://www.mun.ca/biology/scarr/FISH_chromosomes_300dpi.jpg
  - http://agaudi.files.wordpress.com/2008/09/dna_overview_es.png
  - http://joachimj.club.fr/imagesmada2004bis/PlanchePhylogeniedesprimates.jpg
  - http://http://fr.wikipedia.org/wiki/Gene
  - http://www.g-language.org/g3/

# Appendix

1. Appendix
   - Pseudo boolean transformation for other problems
   - ILCS$_X$ and IILCS$_X$
   - Visualization tool
   - Common intervals filtering
   - First experimental results

# Appendix

# Transformation for $ICOM_E$: objective function

Objective:

$$\text{maximize} \sum_{k,l,m,n} I_{k,l,m,n}$$

# Transformation for $ICOM_E$: objective function

Objective:

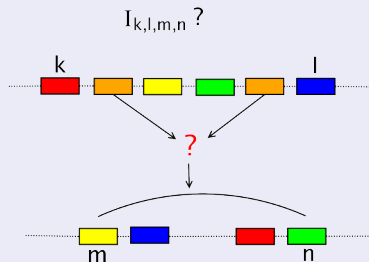$$\text{maximize} \sum_{k,l,m,n} I_{k,l,m,n}$$

Improvements:

- Add rules to decrease the size of the instance

If all orange genes are located
between the red and green one

We must have at least one orange
gene to validate $I_{k,l,m,n}$
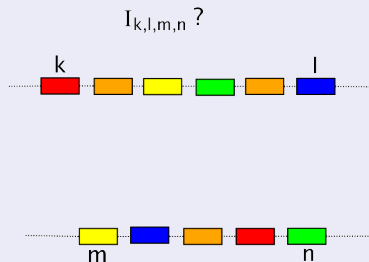
# Transformation for $ICOM_E$: objective function

Objective:

$$\text{maximize} \sum_{k,l,m,n} I_{k,l,m,n}$$

Improvements:

- Add rules to decrease the size of the instance

Else, we do not generate
variable $I_{k,l,m,n}$

# Other problems ?

## Other models

- C1: (Exemplar model)
  $$\forall f \in \mathcal{F}_{G_0} \cup \mathcal{F}_{G_1}, \quad \sum_{\substack{1 \leqslant a \leqslant \eta_{G_0} \\ G_0[a]=f}} \sum_{\substack{1 \leqslant b \leqslant \eta_{G_1} \\ G_1[b]=f}} x_b^a = 1$$
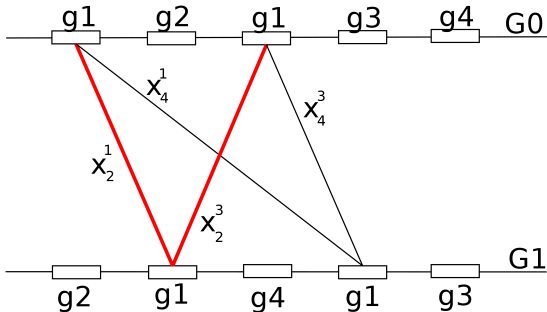
- C1': (Maximal matching model)
  $$\forall f \in \mathcal{F}_{G_0} \cup \mathcal{F}_{G_1}, \quad \sum_{\substack{1 \leqslant a \leqslant \eta_{G_0} \\ G_0[a]=f}} \sum_{\substack{1 \leqslant b \leqslant \eta_{G_1} \\ G_1[b]=f}} x_b^a = \min\{occ(f, G_0), occ(f, G_1)\}$$

- C1": (Intermediate matching model)
  $$\forall f \in \mathcal{F}_{G_0} \cup \mathcal{F}_{G_1}, \quad \sum_{\substack{1 \leqslant a \leqslant \eta_{G_0} \\ G_0[a]=f}} \sum_{\substack{1 \leqslant b \leqslant \eta_{G_1} \\ G_1[b]=f}} x_b^a \geqslant 1$$

# Other models



- $\forall a = 1, 2, \ldots, \eta_{G_0}, \qquad \displaystyle\sum_{\substack{1 \leqslant b \leqslant \eta_{G_1} \\ G_0[a] = G_1[b]}} x_b^a \leqslant 1$

- $\forall b = 1, 2, \ldots, \eta_{G_1}, \qquad \displaystyle\sum_{\substack{1 \leqslant a \leqslant \eta_{G_0} \\ G_0[a] = G_1[b]}} x_b^a \leqslant 1$

# Other problems ?

### Other measures

- $ICONS_X$:
  Generate only variables $I_{k,l,m,n}$ such that
  $$( ( G_0[k] = G_1[m] \wedge G_0[\ell] = G_1[n] ) \vee$$
  $$( G_0[k] = -G_1[n] \wedge G_0[\ell] = -G_1[m] ) )\}$$

# Other problems ?

### Other measures

- $ICONS_X$:
  Generate only variables $I_{k,l,m,n}$ such that
  $(\ (\ G_0[k] = G_1[m] \wedge G_0[\ell] = G_1[n]\ ) \vee$
  $(\ G_0[k] = -G_1[n] \wedge G_0[\ell] = -G_1[m]\ )\ )\}$

- $BD_X$ and $ADJ_X$:
  Other transformation

Angibaud Sébastien                    Defence of Phd Thesis                    October 7th 2009        6 / 23

# Appendix

# ILCS$_M$ heuristic

**LCS:** Longest Common Substring [Tichy, 84]

1 2 3 **4 5 6 7**

**6 7 4 5 1 6** 3 2 1

# ILCS$_M$ heuristic

**LCS:** Longest Common Substring [Tichy, 84]

<div align="center">

**1 2 3** **4 5 6 7**

**6 7 4 5 1 6** **3 2 1**

</div>

ILCS$_M$ heuristic

**Idea:** Match genes of the LCS until saturation

# ILCS$_M$ heuristic

**LCS:** Longest Common Substring [Tichy, 84]

1 2 3 **4 5 6 7**

**6 7 4 5 1 6** 3 2 1

## ILCS$_M$ heuristic

**Idea:** Match genes of the LCS until saturation

1. Compute the Longest Common Substring **S**

# ILCS$_M$ heuristic

**LCS:** Longest Common Substring [Tichy, 84]

$$1\ 2\ 3\ \textbf{4\ 5\ 6\ 7}$$

$$\textbf{6\ 7\ 4\ 5\ 1\ 6}\ 3\ 2\ 1$$

### ILCS$_M$ heuristic

**Idea:** Match genes of the LCS until saturation

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly

# ILCS*M* heuristic

**LCS:** Longest Common Substring [Tichy, 84]

$$1\ 2\ 3\ \textbf{4}\ \textbf{5}\ \textbf{6}\ \textbf{7}$$

$$\textbf{6}\ \textbf{7}\ \textbf{4}\ \textbf{5}\ \textbf{1}\ \textbf{6}\ 3\ 2\ 1$$

## ILCS*M* heuristic

**Idea:** Match genes of the LCS until saturation

1. Compute the Longest Common Substring *S*
2. Match all the genes of *S* accordingly
3. Iterate the process until saturation

# ILCS*M* heuristic

---

**LCS:** Longest Common Substring [Tichy, 84]

<p align="center">1 2 3   4 5   6 7</p>

<p align="center">6 7   4 5   1 6   3 2 1</p>

---

### ILCS*M* heuristic

**Idea:** Match genes of the LCS until saturation

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Iterate the process until saturation

---

# ILCS$_M$ heuristic

**LCS:** Longest Common Substring [Tichy, 84]

1 2 3  4 5  6 7

6 7  4 5  **1 6**  3 2 1

## ILCS$_M$ heuristic

**Idea:** Match genes of the LCS until saturation

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Iterate the process until saturation

# ILCS*M* heuristic

**LCS:** Longest Common Substring [Tichy, 84]

1 2 3  4 5  6 7

6 7  4 5  **1 6**  3 2 1

## ILCS*M* heuristic

**Idea:** Match genes of the LCS until saturation

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Iterate the process until saturation
4. Remove all the genes that have not been matched

Angibaud Sébastien                Defence of Phd Thesis                October 7th 2009      8 / 23

# ILCS*M* heuristic

**LCS:** Longest Common Substring [Tichy, 84]

<div align="center">

1 2 3   4 5   6 7

6 7   4 5   3 2 1

</div>

## ILCS*M* heuristic

**Idea:** Match genes of the LCS until saturation

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Iterate the process until saturation
4. Remove all the genes that have not been matched

# ILCS$_M$ heuristic

**LCS:** Longest Common Substring [Tichy, 84]

1 2 3  4 5  6 7

6 7  4 5  3 2 1

### ILCS$_M$ heuristic

**Idea:** Match genes of the LCS until saturation

1. Compute the Longest Common Substring $S$
2. Match all the genes of $S$ accordingly
3. Iterate the process until saturation
4. Remove all the genes that have not been matched
5. Compute the number of common intervals

# ILCS*M* heuristic

**LCS:** Longest Common Substring [Tichy, 84]

1 2 3    4 5    6 7

6 7    4 5    3 2 1

$\Rightarrow$ number of common intervals $= 19$

## ILCS*M* heuristic

**Idea:** Match genes of the LCS until saturation

1. Compute the Longest Common Substring *S*
2. Match all the genes of *S* accordingly
3. Iterate the process until saturation
4. Remove all the genes that have not been matched
5. Compute the number of common intervals

# IILCS*M* heuristic

**LCS:** Longest Common Substring [Tichy, 84]

1 2 3 **4 5 6 7**

**6 7 4 5 1 6** 3 2 1

### IILCS*M* heuristic

**Idea:** Remove genes that cannot be matched

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Remove genes that cannot be matched

# IILCS*M* heuristic

**LCS:** Longest Common Substring [Tichy, 84]

**1 2 3** **4 5 6 7**

**6 7 4 5** **1** **6** **3 2 1**

### IILCS*M* heuristic

**Idea:** Remove genes that cannot be matched

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Remove genes that cannot be matched

# IILCS*M* heuristic

**LCS:** Longest Common Substring [Tichy, 84]

1 2 3 **4 5 6 7**

**6 7 4 5 6** 3 2 1

## IILCS*M* heuristic

**Idea:** Remove genes that cannot be matched

1. Compute the Longest Common Substring *S*
2. Match all the genes of *S* accordingly
3. Remove genes that cannot be matched

# IILCS$_M$ heuristic

**LCS:** Longest Common Substring [Tichy, 84]

1 2 3 **4 5 6 7**

**6 7 4 5 6** 3 2 1

## IILCS$_M$ heuristic

**Idea:** Remove genes that cannot be matched

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Remove genes that cannot be matched
4. Iterate the process until saturation

# IILCS*M* heuristic

**LCS:** Longest Common Substring [Tichy, 84]

1 2 3  4 5 6  **7**

**6 7**  4 5 6  3 2 1

## IILCS*M* heuristic

**Idea:** Remove genes that cannot be matched

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Remove genes that cannot be matched
4. Iterate the process until saturation

Angibaud Sébastien                     Defence of Phd Thesis                     October 7th 2009     9 / 23

# IILCS*M* heuristic

**LCS:** Longest Common Substring [Tichy, 84]

1 2 3   4 5 6   **7**

**6** **7**   4 5 6   3 2 1

### IILCS*M* heuristic

**Idea:** Remove genes that cannot be matched

1. Compute the Longest Common Substring *S*
2. Match all the genes of *S* accordingly
3. Remove genes that cannot be matched
4. Iterate the process until saturation

# IILCS*M* heuristic

**LCS:** Longest Common Substring [Tichy, 84]

1 2 3   4 5 6   **7**

**7**   4 5 6   3 2 1

## IILCS*M* heuristic

**Idea:** Remove genes that cannot be matched

1. Compute the Longest Common Substring *S*
2. Match all the genes of *S* accordingly
3. Remove genes that cannot be matched
4. Iterate the process until saturation

# IILCS*M* heuristic

**LCS:** Longest Common Substring [Tichy, 84]

$$1\ 2\ 3\quad 4\ 5\ 6\quad 7$$

$$7\quad 4\ 5\ 6\quad 3\ 2\ 1$$

## IILCS*M* heuristic

**Idea:** Remove genes that cannot be matched

1. Compute the Longest Common Substring *S*
2. Match all the genes of *S* accordingly
3. Remove genes that cannot be matched
4. Iterate the process until saturation

# IILCS*M* heuristic

**LCS:** Longest Common Substring [Tichy, 84]

$$1\ 2\ 3\ \ 4\ 5\ 6\ \ 7$$

$$7\ \ 4\ 5\ 6\ \ 3\ 2\ 1$$

## IILCS*M* heuristic

**Idea:** Remove genes that cannot be matched

1. Compute the Longest Common Substring *S*
2. Match all the genes of *S* accordingly
3. Remove genes that cannot be matched
4. Iterate the process until saturation
5. Compute the number of common intervals

Angibaud Sébastien          Defence of Phd Thesis          October 7th 2009          9 / 23

# IILCS$_M$ heuristic

**LCS:** Longest Common Substring [Tichy, 84]

1 2 3  4 5 6  7

7  4 5 6  3 2 1

$\Rightarrow$ number of common intervals = 20

## IILCS$_M$ heuristic

**Idea:** Remove genes that cannot be matched

1. Compute the Longest Common Substring **S**
2. Match all the genes of **S** accordingly
3. Remove genes that cannot be matched
4. Iterate the process until saturation
5. Compute the number of common intervals

# Heuristics: adaptation for other models

### exemplar model

- For each gene family, we keep only the first occurrence in an LCS
- At each iteration, we remove all genes that cannot be matched

# Heuristics: adaptation for other models

### exemplar model

- For each gene family, we keep only the first occurrence in an LCS
- At each iteration, we remove all genes that cannot be matched

### intermediate model

- We stop if, for each gene family, there exists at least one occurrence in the matching

# Experimental results: $ICOM_M$

# Experimental results: $ADJ_E$

# Experimental results: $ADJ_M$

# Experimental results: $ADJ_l$

# Appendix

# Visualization

# Appendix

# Common intervals filtering

- Lots of common intervals
- Relevance of common intervals ?
⇒ Three filters to emphasize *the most interresting* common intervals

# Common intervals filtering

- Lots of common intervals
- Relevance of common intervals ?
⇒ Three filters to emphasize *the most interresting* common intervals

### Filters

1. **Maximal common intervals**:
   Select only common intervals that are not contained in another one

# Common intervals filtering

- Lots of common intervals
- Relevance of common intervals ?
- ⇒ Three filters to emphasize *the most interresting* common intervals

## Filters

**1 Maximal common intervals**:
Select only common intervals that are not contained in another one

**2 Annotated common intervals**:
Select maximal common intervals that contain some annotations in the *Ecocyc database*

# Common intervals filtering

- Lots of common intervals
- Relevance of common intervals ?
⇒ Three filters to emphasize *the most interresting* common intervals

### Filters

1. **Maximal common intervals**:
   Select only common intervals that are not contained in another one

2. **Annotated common intervals**:
   Select maximal common intervals that contain some annotations in the *Ecocyc database*

3. **Relevant common intervals**:
   Select annotated common intervals with good *p-value* (obtained by GO-TermFinder)

# Common intervals filtering

- Lots of common intervals
- Relevance of common intervals ?
⇒ Three filters to emphasize *the most interresting* common intervals

## Filters

1. **Maximal common intervals**:
   Select only common intervals that are not contained in another one

2. **Annotated common intervals**:
   Select maximal common intervals that contain some annotations in the *Ecocyc database*

3. **Relevant common intervals**:
   Select annotated common intervals with good *p-value* (obtained by GO-TermFinder)

# Appendix

# Experimental results

### Input : six chromosomes of $\gamma$-Proteobacteria

| NCBI identifiant | Name |
| --- | --- |
| NC_000913 | *Escherichia coli* K12 |
| NC_002505 | *Vibrio cholerae* 01 biovar eltor str. N16961 chromosome I |
| NC_002506 | *Vibrio cholerae* 01 biovar eltor str. N16961 chromosome II |
| NC_009456 | *Vibrio cholerae* 0395 chromosome I |
| NC_009457 | *Vibrio cholerae* 0395 chromosome II |
| NC_006840 | *Vibrio fischeri* ES114 chromosome I |
| NC_006841 | *Vibrio fischeri* ES114 chromosome II |

# Results: common intervals

| genome $G_2$ | genome size | | method | computational time | | common intervals | |
| | E. coli | $G_2$ | | Inparanoid (s) | matching (s) | number | maximal |
|---|---|---|---|---|---|---|---|
| **NC002505** | **4243** | **2742** | **IILCS** | 1144 | 15 | **7418** | **274** |
| **NC002506** | **4243** | **1093** | **PSB** | 638 | 41 | **246** | **50** |
| **NC009456** | **4243** | **1133** | **PSB** | 651 | 46 | **264** | **55** |
| **NC009457** | **4243** | **2742** | **IILCS** | 1199 | 18 | **7204** | **278** |
| **NC006840** | **4243** | **2586** | **IILCS** | 1012 | 1 | **3865** | **255** |
| **NC006841** | **4243** | **1175** | **IILCS** | 715 | 1 | **203** | **62** |

# Experimental results