


# Computational Aspects in Gene Order Studies

Jens Stoye

AG Genominformatik, Technische Fakultät

Institute of Bioinformatics, Center of Biotechnology

 Bielefeld University, Germany

## Comparative genomics

---

Comparative genomics is applied with several goals in mind:

- Comparative gene finding, annotation
- Binding site prediction, phylogenetic footprinting
- Genome alignment
- Correlated expression
- Conserved gene neighborhood
- Rearrangement studies

## Comparative genomics “at a higher level”

---

Concentrate on large scale layout of the genomes:

- Study genomes based on their *gene order*.
- Represent genomes by their sequence of genes.

## Comparative genomics “at a higher level”

---

Concentrate on large scale layout of the genomes:

- Study genomes based on their *gene order*.
- Represent genomes by their sequence of genes.

genome 1 \_\_\_\_\_

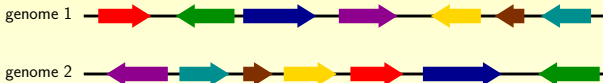
genome 2 \_\_\_\_\_

## Comparative genomics “at a higher level”

---

Concentrate on large scale layout of the genomes:

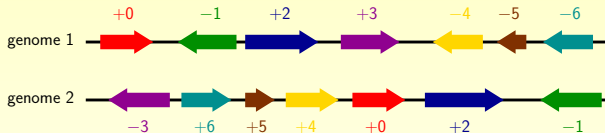
- Study genomes based on their *gene order*.
- Represent genomes by their sequence of genes.



## Comparative genomics “at a higher level”

Concentrate on large scale layout of the genomes:

- Study genomes based on their *gene order*.
- Represent genomes by their sequence of genes.

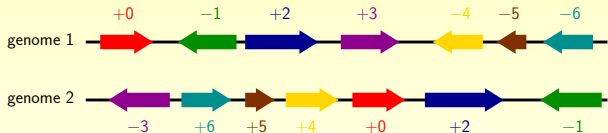


## Comparative genomics “at a higher level”

---

Concentrate on large scale layout of the genomes:

- Study genomes based on their *gene order*.
- Represent genomes by their sequence of genes.



More formally:

- Genes = (signed) elements from the set  $N = \{0, \dots, n\}$ .
- Corresponding (*orthologous*) genes get the same number.
- Genomes = (signed) permutations of  $N$ .

## **Protein function prediction**

---

The yeast *S. cerevisiae* was sequenced in 1995.

Still, about 30% of the ORFs in the *MIPS Yeast Database* have no functional annotation.

Functional annotation is time consuming and expensive.



## **Protein function prediction**

---

The yeast *S. cerevisiae* was sequenced in 1995.

Still, about 30% of the ORFs in the *MIPS Yeast Database* have no functional annotation.

Functional annotation is time consuming and expensive.

- In the lab:
  - Genetical and biochemical analysis
  - Correlated expression

## **Protein function prediction**

---

The yeast *S. cerevisiae* was sequenced in 1995.

Still, about 30% of the ORFs in the *MIPS Yeast Database* have no functional annotation.

Functional annotation is time consuming and expensive.

- In the lab:
  - Genetical and biochemical analysis
  - Correlated expression
- Homology based:
  - Protein families
  - Functional domains

## Protein function prediction

---

The yeast *S. cerevisiae* was sequenced in 1995.

Still, about 30% of the ORFs in the *MIPS Yeast Database* have no functional annotation.

Functional annotation is time consuming and expensive.

- In the lab:
  - Genetical and biochemical analysis
  - Correlated expression
- Homology based:
  - Protein families
  - Functional domains
- Genome based:
  - Rosetta stone method (gene fusion, domain fusion)
  - Phylogenetic profiles (correlated evolution)
  - Gene order (co-occurrence of genes in genomes)

## Protein function prediction

---

The yeast *S. cerevisiae* was sequenced in 1995.

Still, about 30% of the ORFs in the *MIPS Yeast Database* have no functional annotation.

Functional annotation is time consuming and expensive.

- In the lab:
  - Genetical and biochemical analysis
  - Correlated expression
- Homology based:
  - Protein families
  - Functional domains
- Genome based:
  - Rosetta stone method (gene fusion, domain fusion)
  - Phylogenetic profiles (correlated evolution)
  - Gene order (co-occurrence of genes in genomes)
- Literature based:
  - Natural language processing

## Genome-based gene function prediction

---

*Comparative genomics meets functional genomics.*

**Idea:** Genes that repeatedly cluster together in phylogenetically remotely related genomes are functionally associated:

- interacting proteins
- proteins of the same protein complex
- enzymes of the same metabolic pathway

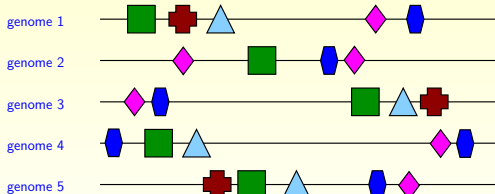
## Genome-based gene function prediction

---

*Comparative genomics meets functional genomics.*

**Idea:** Genes that repeatedly cluster together in phylogenetically remotely related genomes are functionally associated:

- interacting proteins
- proteins of the same protein complex
- enzymes of the same metabolic pathway



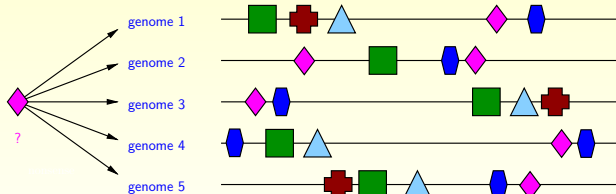
## Genome-based gene function prediction

---

*Comparative genomics meets functional genomics.*

**Idea:** Genes that repeatedly cluster together in phylogenetically remotely related genomes are functionally associated:

- interacting proteins
- proteins of the same protein complex
- enzymes of the same metabolic pathway



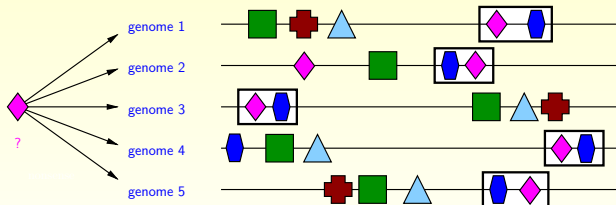
## Genome-based gene function prediction

---

*Comparative genomics meets functional genomics.*

**Idea:** Genes that repeatedly cluster together in phylogenetically remotely related genomes are functionally associated:

- interacting proteins
- proteins of the same protein complex
- enzymes of the same metabolic pathway





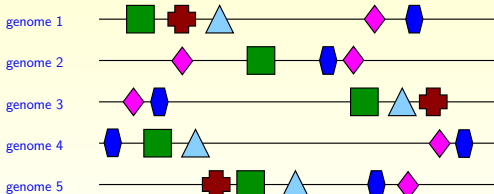
## Genome-based gene function prediction

---

*Comparative genomics meets functional genomics.*

**Idea:** Genes that repeatedly cluster together in phylogenetically remotely related genomes are functionally associated:

- interacting proteins
- proteins of the same protein complex
- enzymes of the same metabolic pathway



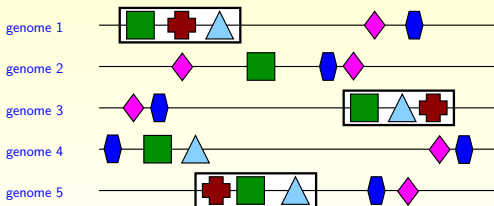
## Genome-based gene function prediction

---

*Comparative genomics meets functional genomics.*

**Idea:** Genes that repeatedly cluster together in phylogenetically remotely related genomes are functionally associated:

- interacting proteins
- proteins of the same protein complex
- enzymes of the same metabolic pathway



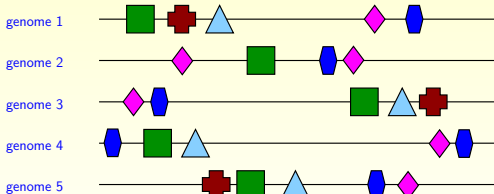
## Genome-based gene function prediction

---

*Comparative genomics meets functional genomics.*

**Idea:** Genes that repeatedly cluster together in phylogenetically remotely related genomes are functionally associated:

- interacting proteins
- proteins of the same protein complex
- enzymes of the same metabolic pathway



## Genome-based gene function prediction

---

*Comparative genomics meets functional genomics.*

**Idea:** Genes that repeatedly cluster together in phylogenetically remotely related genomes are functionally associated:

- interacting proteins
- proteins of the same protein complex
- enzymes of the same metabolic pathway

## Genome-based gene function prediction

---

*Comparative genomics meets functional genomics.*

**Idea:** Genes that repeatedly cluster together in phylogenetically remotely related genomes are functionally associated:

- interacting proteins
- proteins of the same protein complex
- enzymes of the same metabolic pathway

Marcotte *et al.*: Detecting protein function and protein-protein interactions from genome sequences. *Science*, 1999.

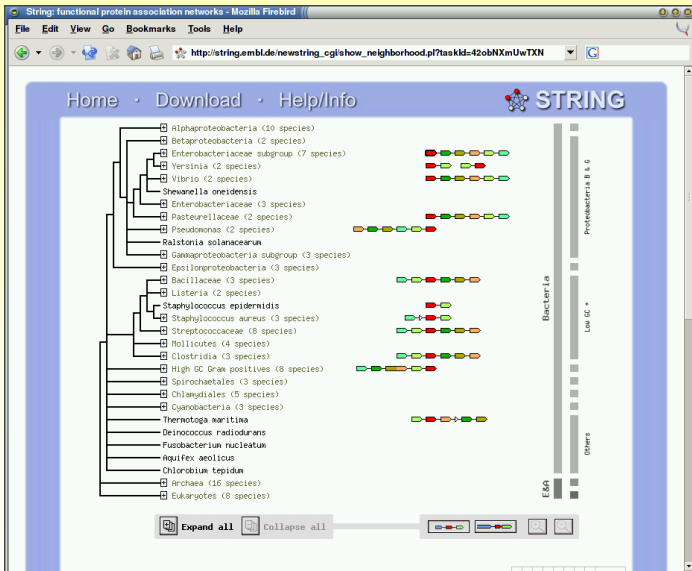
Overbeek *et al.*: The use of gene clusters to infer functional coupling. *PNAS*, 1999.

Snel *et al.*: STRING: A web-server to retrieve and display the repeatedly occurring neighbourhood of a gene, *NAR*, 2000.

Zheng *et al.*: Phylogenetic detection of conserved gene clusters in microbial genomes, *BMC Bioinformatics*, 2005.

# STRING Web server (Snel et al., 2000)

<http://string.embl.de/>



# STRING Web server (Snel et al., 2000)

<http://string.embl.de/>

String: functional protein association networks - Mozilla Firebird

File Edit View Go Bookmarks Tools Help

[http://string.embl.de/newstring.cgi/show\\_neighborhood.pl?tas\\_kkd=42obNXmUwTXN](http://string.embl.de/newstring.cgi/show_neighborhood.pl?tas_kkd=42obNXmUwTXN)

Phylogenetic tree showing relationships between various bacterial species and groups, including Streptococcaceae, Mollicutes, Clostridia, High GC Gram positives, Spirochaetales, Chlamydiales, Cyanobacteria, Thermotoga maritima, Deinococcus radiodurans, Fusobacterium nucleatum, Aquifex aeolicus, Chlorobium tepidum, Archaea, and Eukaryotes.

Expand all Collapse all

**Your Input:**

- RBS\_D\_ECOLI High affinity ribose transport protein rbsD (151 aa)

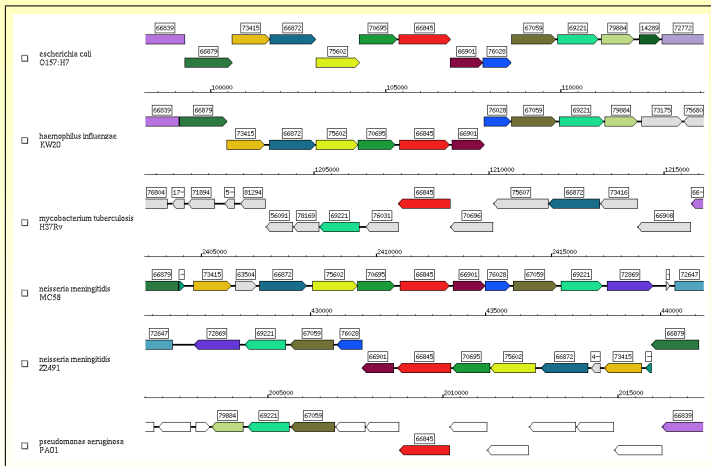
**Predicted Functional Associations:**

	Neighborhood	Gene Fusion	Cooccurrence	[Homo]logJ	Coexpression	Experiments	Databases	Textmining	Score
<input checked="" type="checkbox"/> RBSB_ECOLI D-ribose-binding periplasmic protein precursor (296 aa)	●	●	●	●	●	●	●	●	0.973
<input checked="" type="checkbox"/> RBS_C_ECOLI Ribose transport system permease protein rbsC (321 aa)	●	●	●	●	●	●	●	●	0.966
<input checked="" type="checkbox"/> RBSK_ECOLI Ribokinase (EC 2.7.1.15) (309 aa)	●	●	●	●	●	●	●	●	0.957
<input checked="" type="checkbox"/> RBSR_ECOLI Ribose transport ATP-binding protein rbsK (501 aa)	●	●	●	●	●	●	●	●	0.910
<input checked="" type="checkbox"/> RBSL_ECOLI Ribose operon repressor (329 aa)	●	●	●	●	●	●	●	●	0.832
<input checked="" type="checkbox"/> FUCU_ECOLI Fucose operon fucU protein (140 aa)	●	●	●	●	●	●	●	●	0.631
<input checked="" type="checkbox"/> ALS_C_ECOLI D-allose transport system permease protein alsC (326 aa)	●	●	●	●	●	●	●	●	0.592
<input checked="" type="checkbox"/> YPHD_ECOLI Hypothetical ABC transporter permease protein yphD (332 aa)	●	●	●	●	●	●	●	●	0.564
<input checked="" type="checkbox"/> ALSB_ECOLI D-allose-binding periplasmic protein precursor (ALBP) (311 aa)	●	●	●	●	●	●	●	●	0.494
<input checked="" type="checkbox"/> UWAC_ECOLI Uronate isomerase (EC 5.3.1.12) (Glucuronate isomerase) (Uronic isomer [...])	●	●	●	●	●	●	●	●	0.430

**Views:**

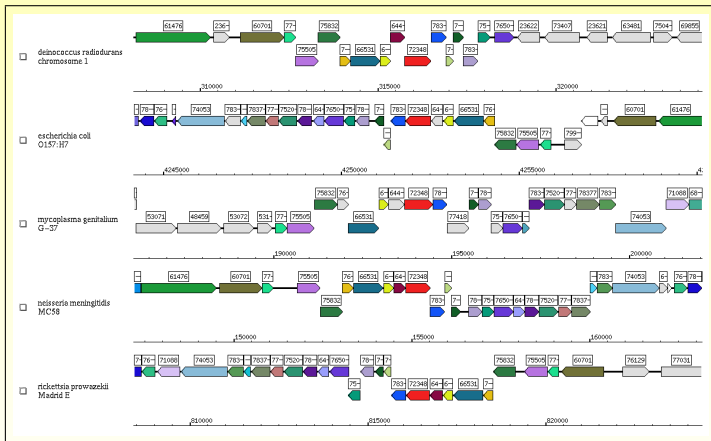
- Neighborhood
- Fusion
- Cooccurrence
- Coexpression
- Experiments
- Databases
- Textmining
- Summary Network

# Genome Windows: DCW cluster (division and cell wall)



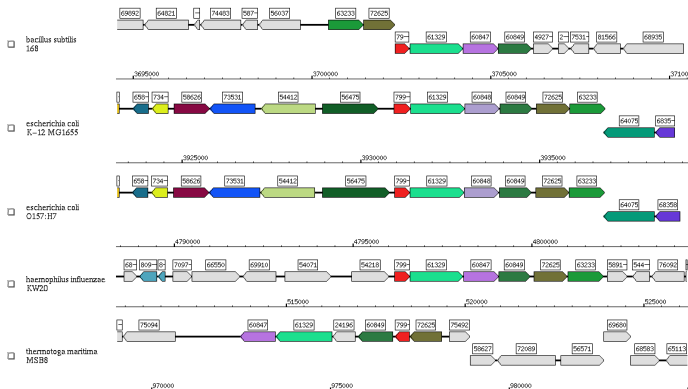


# Genome Windows: Ribosomal Super Operon



# Genome Windows: Ribose ABC Transporter

## Genome Windows of Cluster 79903:



## Formalization of gene clusters: common intervals

---

Given permutations (genomes)  $\pi_1, \pi_2, \dots, \pi_k$  of the numbers (genes)  $0, 1, \dots, n$ , find subsets of numbers that occur **contiguously in all permutations.**

$\pi_1$	0	1	2	3	4	5	6	7
$\pi_2$	6	7	5	1	4	3	2	0
$\pi_3$	7	0	1	2	4	3	5	6

## Formalization of gene clusters: common intervals

---

Given permutations (genomes)  $\pi_1, \pi_2, \dots, \pi_k$  of the numbers (genes)  $0, 1, \dots, n$ , find subsets of numbers that occur contiguously in all permutations.

$\pi_1$	0	1	2	3	4	5	6	7
$\pi_2$	6	7	5	1	4	3	2	0
$\pi_3$	7	0	1	2	4	3	5	6

Common intervals: [3,4]

## Formalization of gene clusters: common intervals

---

Given permutations (genomes)  $\pi_1, \pi_2, \dots, \pi_k$  of the numbers (genes)  $0, 1, \dots, n$ , find subsets of numbers that occur contiguously in all permutations.

$\pi_1$	0	1	2	3	4	5	6	7
$\pi_2$	6	7	5	1	4	3	2	0
$\pi_3$	7	0	1	2	4	3	5	6

Common intervals: [3,4] [2,4]

## Formalization of gene clusters: common intervals

---

Given permutations (genomes)  $\pi_1, \pi_2, \dots, \pi_k$  of the numbers (genes)  $0, 1, \dots, n$ , find subsets of numbers that occur contiguously in all permutations.

$\pi_1$	0	1	2	3	4	5	6	7
$\pi_2$	6	7	5	1	4	3	2	0
$\pi_3$	7	0	1	2	4	3	5	6

Common intervals: [3,4] [2,4] [1,4]

## Formalization of gene clusters: common intervals

---

Given permutations (genomes)  $\pi_1, \pi_2, \dots, \pi_k$  of the numbers (genes)  $0, 1, \dots, n$ , find subsets of numbers that occur contiguously in all permutations.

$\pi_1$	0 1 2 3 4 5 6 7
$\pi_2$	6 7 5 1 4 3 2 0
$\pi_3$	7 0 1 2 4 3 5 6

Common intervals: [3,4] [2,4] [1,4] [0,4]

## Formalization of gene clusters: common intervals

---

Given permutations (genomes)  $\pi_1, \pi_2, \dots, \pi_k$  of the numbers (genes)  $0, 1, \dots, n$ , find subsets of numbers that occur contiguously in all permutations.

$\pi_1$	0	1	2	3	4	5	6	7
$\pi_2$	6	7	5	1	4	3	2	0
$\pi_3$	7	0	1	2	4	3	5	6

Common intervals: [3,4] [2,4] [1,4] [0,4] [1,5]



## Formalization of gene clusters: common intervals

---

Given permutations (genomes)  $\pi_1, \pi_2, \dots, \pi_k$  of the numbers (genes)  $0, 1, \dots, n$ , find subsets of numbers that occur contiguously in all permutations.

$\pi_1$	0 1 2 3 4 5 6 7
$\pi_2$	6 7 5 1 4 3 2 0
$\pi_3$	7 0 1 2 4 3 5 6

Common intervals: [3,4] [2,4] [1,4] [0,4] [1,5] [0,5]

## Formalization of gene clusters: common intervals

---

Given permutations (genomes)  $\pi_1, \pi_2, \dots, \pi_k$  of the numbers (genes)  $0, 1, \dots, n$ , find subsets of numbers that occur contiguously in all permutations.

$\pi_1$	0 1 2 3 4 5 6 7
$\pi_2$	6 7 5 1 4 3 2 0
$\pi_3$	7 0 1 2 4 3 5 6

Common intervals: [3,4] [2,4] [1,4] [0,4] [1,5] [0,5] [0,7]

## Formalization of gene clusters: common intervals

---

Given permutations (genomes)  $\pi_1, \pi_2, \dots, \pi_k$  of the numbers (genes)  $0, 1, \dots, n$ , find subsets of numbers that occur **contiguously in all permutations**.

$\pi_1$	0	1	2	3	4	5	6	7
$\pi_2$	6	7	5	1	4	3	2	0
$\pi_3$	7	0	1	2	4	3	5	6

Common intervals: [3,4] [2,4] [1,4] [0,4] [1,5] [0,5] [0,7]

### Algorithms:

- Uno & Yagiura, *Algorithmica* 2000:  
Find all common intervals of 2 permutations in  $\mathcal{O}(n + |\text{output}|)$  time.
- Heber & JS, *CPM* 2001:  
Find all common intervals of  $k \geq 2$  permutations in  $\mathcal{O}(kn + |\text{output}|)$  time.

## Finding all common intervals of two permutations $\pi_1$ and $\pi_2$

Let  $1 \leq x \leq y \leq n$ .

**Notation:**  $\pi([x, y]) := \{\pi(x), \pi(x+1), \dots, \pi(y)\}$

**Definitions:**

$$l(x, y) := \min \pi_2([x, y])$$
$$u(x, y) := \max \pi_2([x, y])$$
$$f(x, y) := u(x, y) - l(x, y) - (y - x)$$

**Example:**

$\pi_1$	0	1	2	3	4	5	6	7
	0	1	2	3	4	5	6	7
$\pi_2$	6	7	5	1	4	3	2	0
	0	1	2	3	4	5	6	7

## Finding all common intervals of two permutations $\pi_1$ and $\pi_2$

Let  $1 \leq x \leq y \leq n$ .

**Notation:**  $\pi([x, y]) := \{\pi(x), \pi(x+1), \dots, \pi(y)\}$

**Definitions:**  
 $l(x, y) := \min \pi_2([x, y])$   
 $u(x, y) := \max \pi_2([x, y])$   
 $f(x, y) := u(x, y) - l(x, y) - (y - x)$

**Example:**

$\pi_1$	<table border="1"><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td></tr><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td></tr></table>	0	1	2	3	4	5	6	7	0	1	2	3	4	5	6	7	$f(3, 6) = 4 - 1 - (6 - 3) = 0$
0	1	2	3	4	5	6	7											
0	1	2	3	4	5	6	7											
$\pi_2$	<table border="1"><tr><td>6</td><td>7</td><td>5</td><td>1</td><td>4</td><td>3</td><td>2</td><td>0</td></tr><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td></tr></table>	6	7	5	1	4	3	2	0	0	1	2	3	4	5	6	7	
6	7	5	1	4	3	2	0											
0	1	2	3	4	5	6	7											

## Finding all common intervals of two permutations $\pi_1$ and $\pi_2$

Let  $1 \leq x \leq y \leq n$ .

**Notation:**  $\pi([x, y]) := \{\pi(x), \pi(x+1), \dots, \pi(y)\}$

**Definitions:**  $l(x, y) := \min \pi_2([x, y])$   
 $u(x, y) := \max \pi_2([x, y])$   
 $f(x, y) := u(x, y) - l(x, y) - (y - x)$

**Example:**

$\pi_1$	<table border="1"><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td></tr><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td></tr></table>	0	1	2	3	4	5	6	7	0	1	2	3	4	5	6	7	$f(3, 6) = 4 - 1 - (6 - 3) = 0$
0	1	2	3	4	5	6	7											
0	1	2	3	4	5	6	7											
$\pi_2$	<table border="1"><tr><td>6</td><td>7</td><td>5</td><td>1</td><td>4</td><td>3</td><td>2</td><td>0</td></tr><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td></tr></table>	6	7	5	1	4	3	2	0	0	1	2	3	4	5	6	7	$f(1, 4) = 7 - 1 - (4 - 1) = 3 > 0$
6	7	5	1	4	3	2	0											
0	1	2	3	4	5	6	7											

## Finding all common intervals of two permutations $\pi_1$ and $\pi_2$

Let  $1 \leq x \leq y \leq n$ .

**Notation:**  $\pi([x, y]) := \{\pi(x), \pi(x+1), \dots, \pi(y)\}$

**Definitions:**  
 $l(x, y) := \min \pi_2([x, y])$   
 $u(x, y) := \max \pi_2([x, y])$   
 $f(x, y) := u(x, y) - l(x, y) - (y - x)$

**Example:**

$\pi_1$	<table border="1"><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td></tr><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td></tr></table>	0	1	2	3	4	5	6	7	0	1	2	3	4	5	6	7	$f(3, 6) = 4 - 1 - (6 - 3) = 0$
0	1	2	3	4	5	6	7											
0	1	2	3	4	5	6	7											
$\pi_2$	<table border="1"><tr><td>6</td><td>7</td><td>5</td><td>1</td><td>4</td><td>3</td><td>2</td><td>0</td></tr><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td></tr></table>	6	7	5	1	4	3	2	0	0	1	2	3	4	5	6	7	$f(1, 4) = 7 - 1 - (4 - 1) = 3 > 0$
6	7	5	1	4	3	2	0											
0	1	2	3	4	5	6	7											

**Simple algorithm:** For all  $1 \leq x \leq y \leq n$  test if  $f(x, y) = 0$ .

## Finding all common intervals of two permutations $\pi_1$ and $\pi_2$

Let  $1 \leq x \leq y \leq n$ .

Notation:  $\pi([x, y]) := \{\pi(x), \pi(x+1), \dots, \pi(y)\}$

Definitions:  $l(x, y) := \min \pi_2([x, y])$   
 $u(x, y) := \max \pi_2([x, y])$   
 $f(x, y) := u(x, y) - l(x, y) - (y - x)$

Example:

$\pi_1$	<table border="1"><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td></tr><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td></tr></table>	0	1	2	3	4	5	6	7	0	1	2	3	4	5	6	7	$f(3, 6) = 4 - 1 - (6 - 3) = 0$
0	1	2	3	4	5	6	7											
0	1	2	3	4	5	6	7											
$\pi_2$	<table border="1"><tr><td>6</td><td>7</td><td>5</td><td>1</td><td>4</td><td>3</td><td>2</td><td>0</td></tr><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td></tr></table>	6	7	5	1	4	3	2	0	0	1	2	3	4	5	6	7	$f(1, 4) = 7 - 1 - (4 - 1) = 3 > 0$
6	7	5	1	4	3	2	0											
0	1	2	3	4	5	6	7											

Simple algorithm: For all  $1 \leq x \leq y \leq n$  test if  $f(x, y) = 0$ .

Analysis:  $\mathcal{O}(n^2)$  time.



## Finding all common intervals of two permutations $\pi_1$ and $\pi_2$

Uno & Yagiura, 2000:

Perform the test  $f(x, y) = 0$  not for all pairs  $(x, y)$ .

**Definition:**

For given  $x$ , call a value of  $y > x$  *wasteful*, if and only if for all  $x' \leq x$ :

$$f(x', y) > 0.$$

**Lemma:**

For fixed  $x$ ,  $f(x, y)$  increases monotonically for the non-wasteful indices  $y (> x)$ .

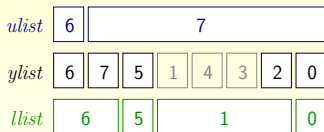
**Algorithm (Idea):**

- $x$  runs in right-to-left direction through a doubly linked list *ylist* that initially contains the entries of  $\pi_2$ .
- In each step, the entries of wasteful indices  $y (> x)$  are removed.
- Test for the remaining  $y > x$  in *ylist* from left to right if  $f(x, y) = 0$ .

## Algorithm RC (Uno & Yagiura)

---

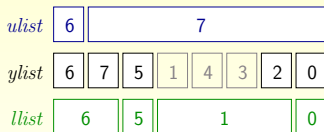
- Removal of wasteful indices from *ylist* is done by means of two additional lists *l*list and *u*list that implement the functions *l* and *u*.
- The elements of *l*list and *u*list are maximal intervals of *ylist* with the same smallest resp. largest element.



## Algorithm RC (Uno & Yagiura)

---

- Removal of wasteful indices from *ylist* is done by means of two additional lists *l*list and *u*list that implement the functions *l* and *u*.
- The elements of *l*list and *u*list are maximal intervals of *ylist* with the same smallest resp. largest element.



Analysis:

$\mathcal{O}(n + |\text{output}|)$  time,  $\mathcal{O}(n)$  space.

## Finding all common intervals of $k \geq 2$ permutations

---

Obvious generalization:

Given  $k$  permutations  $\pi_1, \pi_2, \dots, \pi_k$ .

For  $j = 2, 3, \dots, k$  compute the common intervals of  $\pi_1$  and  $\pi_j$ .

Output all intervals that are found in all of these comparisons.

$\pi_1$	0	1	2	3	4	5	6	7
$\pi_2$	6	7	5	1	4	3	2	0
$\pi_3$	7	0	1	2	4	3	5	6

## Finding all common intervals of $k \geq 2$ permutations

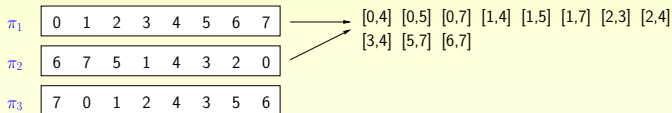
---

Obvious generalization:

Given  $k$  permutations  $\pi_1, \pi_2, \dots, \pi_k$ .

For  $j = 2, 3, \dots, k$  compute the common intervals of  $\pi_1$  and  $\pi_j$ .

Output all intervals that are found in all of these comparisons.



## Finding all common intervals of $k \geq 2$ permutations

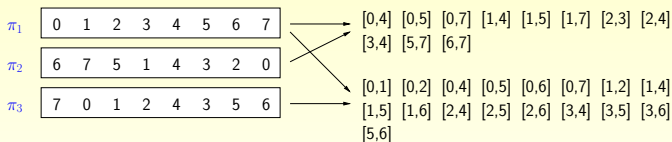
---

Obvious generalization:

Given  $k$  permutations  $\pi_1, \pi_2, \dots, \pi_k$ .

For  $j = 2, 3, \dots, k$  compute the common intervals of  $\pi_1$  and  $\pi_j$ .

Output all intervals that are found in all of these comparisons.



## Finding all common intervals of $k \geq 2$ permutations

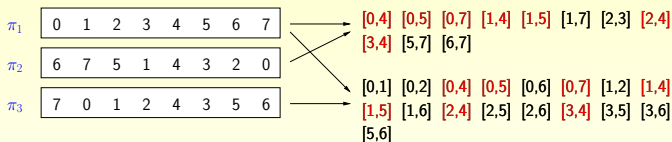
---

Obvious generalization:

Given  $k$  permutations  $\pi_1, \pi_2, \dots, \pi_k$ .

For  $j = 2, 3, \dots, k$  compute the common intervals of  $\pi_1$  and  $\pi_j$ .

Output all intervals that are found in all of these comparisons.



## Finding all common intervals of $k \geq 2$ permutations

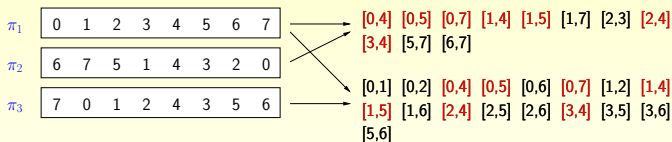
---

Obvious generalization:

Given  $k$  permutations  $\pi_1, \pi_2, \dots, \pi_k$ .

For  $j = 2, 3, \dots, k$  compute the common intervals of  $\pi_1$  and  $\pi_j$ .

Output all intervals that are found in all of these comparisons.



Analysis:

$\mathcal{O}(kn + \sum |K_i|)$  time

where  $K_i$  = the number of common intervals of  $\pi_1$  and  $\pi_i$ .

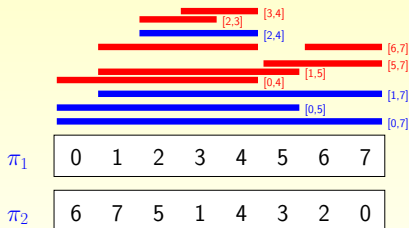


## Irreducible Intervals

---

**Goal:** An algorithm with output-dependent time complexity  $\mathcal{O}(kn + |\text{output}|)$ .

**Observation:** Common intervals form “chains” of non-trivially overlapping intervals.

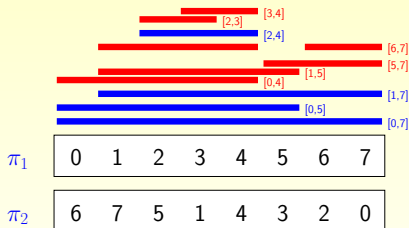


## Irreducible Intervals

---

**Goal:** An algorithm with output-dependent time complexity  $\mathcal{O}(kn + |\text{output}|)$ .

**Observation:** Common intervals form “chains” of non-trivially overlapping intervals.



**Definition:**

A common interval  $c$  is **reducible** if there exists a non-trivial chain that generates  $c$ , otherwise it is **irreducible**.

# Properties of irreducible intervals

---

## Lemma:

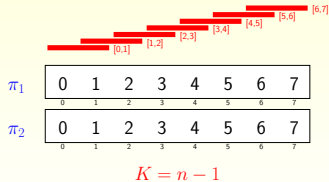
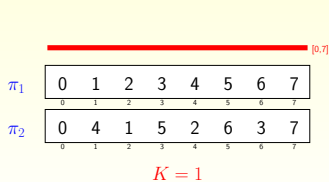
The subchains of all the maximal chains of irreducible intervals generate exactly all common intervals.

## Theorem:

For the number of irreducible intervals  $K$  the following holds:

$$1 \leq K \leq n - 1$$

## Example:

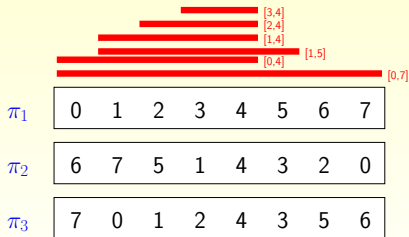


## Finding all common intervals of $k \geq 2$ permutations

---

### Algorithm:

- Find the set of all irreducible intervals.
- Partition this set into maximal chains of non-trivially overlapping intervals.
- For each such chain generate all subchains: the common intervals.

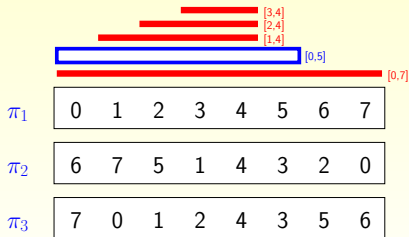


## Finding all common intervals of $k \geq 2$ permutations

---

### Algorithm:

- Find the set of all irreducible intervals.
- Partition this set into maximal chains of non-trivially overlapping intervals.
- For each such chain generate all subchains: the common intervals.

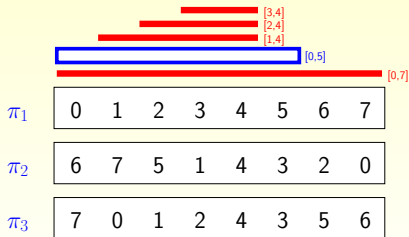


## Finding all common intervals of $k \geq 2$ permutations

---

### Algorithm:

- Find the set of all irreducible intervals.
- Partition this set into maximal chains of non-trivially overlapping intervals.
- For each such chain generate all subchains: the common intervals.



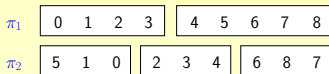
Analysis:  $\mathcal{O}(kn + |\text{output}|)$  time,  $\mathcal{O}(n)$  additional space

## More realistic genome models

---

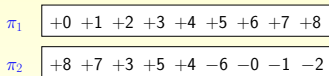
- ① Genomes of higher organisms often have more than one chromosome

⇒ *multichromosomal permutations*



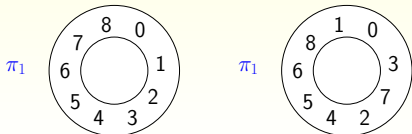
- ② Genes of a cluster should lie on the same DNA strand

⇒ *signed permutations*



- ③ Bacterial, archaeal, and mitochondrial DNA is often circular

⇒ *circular permutations*



## Further extensions of the model

---

### Generalization

- Allow paralogous genes (gene families)

### Relaxation:

- Common intervals do not have to be contained in all genomes
- Not each gene has always to occur

### Other applications

- Inverse of *Consecutive Arrangement Problem*
- Candidates for *Subtour Exchange Crossover* in genetic algorithms
- Model-free distance for phylogenetic analyses



## Inclusion of paralogous genes

---

### Problem:

In case of *duplicated genes*, it is difficult to assign correct orthologous gene pairs.

Possibly *the* ortholog does not even exist.

### Consequence:

Do not distinguish between paralogous gene copies.

### New model:

Use the same element (number) more than once for paralogous copies of genes.

→ genomes are modeled as *sequences* instead of permutations.

## Formal model

---

**Given:**  $k$  sequences  $\mathcal{S} = (S_1, S_2, \dots, S_k)$  over an alphabet  $\Sigma$ .

**Common interval:**

a subset  $C \subseteq \Sigma$  whose elements occur contiguously in each  $S_l \in \mathcal{S}$ .

**Goal:**

Find all maximal occurrences of common intervals in  $\mathcal{S}$ .

**Example:**

$S_1$	3 1 2 3 1 5 2 6
$S_2$	4 3 5 5 5 1 4 2 2
$S_3$	7 5 1 5 3 6 5

## Formal model

---

**Given:**  $k$  sequences  $\mathcal{S} = (S_1, S_2, \dots, S_k)$  over an alphabet  $\Sigma$ .

**Common interval:**

a subset  $C \subseteq \Sigma$  whose elements occur contiguously in each  $S_l \in \mathcal{S}$ .

**Goal:**

Find all maximal occurrences of common intervals in  $\mathcal{S}$ .

**Example:**

$S_1$	3	1	2	3	1	5	2	6	
$S_2$	4	3	5	5	5	1	4	2	2
$S_3$	7	5	1	5	3	6	5		

**Common intervals:**  $\{1,3,5\}$

## Formal model

---

**Given:**  $k$  sequences  $\mathcal{S} = (S_1, S_2, \dots, S_k)$  over an alphabet  $\Sigma$ .

**Common interval:**

a subset  $C \subseteq \Sigma$  whose elements occur contiguously in each  $S_l \in \mathcal{S}$ .

**Goal:**

Find all maximal occurrences of common intervals in  $\mathcal{S}$ .

**Example:**

$S_1$	3	1	2	3	1	5	2	6	
$S_2$	4	3	5	5	5	1	4	2	2
$S_3$	7	5	1	5	3	6	5		

Common intervals:  $\{1,3,5\}$   $\{1,5\}$

## Formal model

---

**Given:**  $k$  sequences  $\mathcal{S} = (S_1, S_2, \dots, S_k)$  over an alphabet  $\Sigma$ .

**Common interval:**

a subset  $C \subseteq \Sigma$  whose elements occur contiguously in each  $S_l \in \mathcal{S}$ .

**Goal:**

Find all maximal occurrences of common intervals in  $\mathcal{S}$ .

**Example:**

$S_1$	3	1	2	3	1	5	2	6	
$S_2$	4	3	5	5	5	1	4	2	2
$S_3$	7	5	1	5	3	6	5		

Common intervals:  $\{1,3,5\}$   $\{1,5\}$   $\{5\}$

## Formal model

---

**Given:**  $k$  sequences  $\mathcal{S} = (S_1, S_2, \dots, S_k)$  over an alphabet  $\Sigma$ .

**Common interval:**

a subset  $C \subseteq \Sigma$  whose elements occur contiguously in each  $S_l \in \mathcal{S}$ .

**Goal:**

Find all maximal occurrences of common intervals in  $\mathcal{S}$ .

**Example:**

$S_1$	<table border="1"><tr><td>3</td><td>1</td><td>2</td><td>3</td><td>1</td><td>5</td><td>2</td><td>6</td></tr></table>	3	1	2	3	1	5	2	6	
3	1	2	3	1	5	2	6			
$S_2$	<table border="1"><tr><td>4</td><td>3</td><td>5</td><td>5</td><td>5</td><td>1</td><td>4</td><td>2</td><td>2</td></tr></table>	4	3	5	5	5	1	4	2	2
4	3	5	5	5	1	4	2	2		
$S_3$	<table border="1"><tr><td>7</td><td>5</td><td>1</td><td>5</td><td>3</td><td>6</td><td>5</td></tr></table>	7	5	1	5	3	6	5		
7	5	1	5	3	6	5				

**Common intervals:**  $\{1,3,5\}$   $\{1,5\}$   $\{5\}$   $\{3\}$

## Formal model

---

**Given:**  $k$  sequences  $\mathcal{S} = (S_1, S_2, \dots, S_k)$  over an alphabet  $\Sigma$ .

**Common interval:**

a subset  $C \subseteq \Sigma$  whose elements occur contiguously in each  $S_l \in \mathcal{S}$ .

**Goal:**

Find all maximal occurrences of common intervals in  $\mathcal{S}$ .

**Example:**

$S_1$	3	1	2	3	1	5	2	6	
$S_2$	4	3	5	5	5	1	4	2	2
$S_3$	7	5	1	5	3	6	5		

Common intervals:  $\{1,3,5\}$   $\{1,5\}$   $\{5\}$   $\{3\}$   $\{1\}$

## Formal model

---

**Given:**  $k$  sequences  $\mathcal{S} = (S_1, S_2, \dots, S_k)$  over an alphabet  $\Sigma$ .

**Common interval:**

a subset  $C \subseteq \Sigma$  whose elements occur contiguously in each  $S_l \in \mathcal{S}$ .

**Goal:**

Find all maximal occurrences of common intervals in  $\mathcal{S}$ .

**Example:**

$S_1$	3 1 2 3 1 5 2 6
$S_2$	4 3 5 5 5 1 4 2 2
$S_3$	7 5 1 5 3 6 5

Common intervals:  $\{1,3,5\}$   $\{1,5\}$   $\{5\}$   $\{3\}$   $\{1\}$



## An elementary algorithm for two sequences

---

Preprocessing: compute two tables for  $S_1 = (3, 1, 2, 3, 1, 5, 2, 6)$ :

$POS[1] = 2, 5$   
 $POS[2] = 3, 7$   
 $POS[3] = 1, 4$   
 $POS[4] = \text{empty}$   
 $POS[5] = 6$   
 $POS[6] = 8$

$NUM(i, j)$ :

$i \setminus j$	0	1	2	3	4	5	6	7
0	1	2	3	3	3	4	4	5
1		1	2	3	3	4	4	5
2			1	2	3	4	4	5
3				1	2	3	4	5
4					1	2	3	4
5						1	2	3
6							1	2
7								1

*Connecting Intervals Algorithm* (Schmidt & JS, CPM 2004):

While reading  $S_2$ , mark in  $S_1$  the observed characters and track maximal intervals of marked characters.

$S_1$ 

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

$S_2$ 

4	3	5	5	5	1	4	2	2
---	---	---	---	---	---	---	---	---

## An elementary algorithm for two sequences

Preprocessing: compute two tables for  $S_1 = (3, 1, 2, 3, 1, 5, 2, 6)$ :

$POS[1] = 2, 5$   
 $POS[2] = 3, 7$   
 $POS[3] = 1, 4$   
 $POS[4] = \text{empty}$   
 $POS[5] = 6$   
 $POS[6] = 8$

$NUM(i, j)$ :

$i \setminus j$	0	1	2	3	4	5	6	7
0	1	2	3	3	3	4	4	5
1		1	2	3	3	4	4	5
2			1	2	3	4	4	5
3				1	2	3	4	5
4					1	2	3	4
5						1	2	3
6							1	2
7								1

*Connecting Intervals Algorithm* (Schmidt & JS, CPM 2004):

While reading  $S_2$ , mark in  $S_1$  the observed characters and track maximal intervals of marked characters.

$S_1$ 

<u>3</u>	1	2	<u>3</u>	1	5	2	6
0	1	2	3	4	5	6	7

$S_2$ 

4	3	5	5	5	1	4	2	2
---	---	---	---	---	---	---	---	---

  
 $i = j$

## An elementary algorithm for two sequences

Preprocessing: compute two tables for  $S_1 = (3, 1, 2, 3, 1, 5, 2, 6)$ :

$POS[1] = 2, 5$   
 $POS[2] = 3, 7$   
 $POS[3] = 1, 4$   
 $POS[4] = \text{empty}$   
 $POS[5] = 6$   
 $POS[6] = 8$

$NUM(i, j)$ :

$i \setminus j$	0	1	2	3	4	5	6	7
0	1	2	3	3	3	4	4	5
1		1	2	3	3	4	4	5
2			1	2	3	4	4	5
3				1	2	3	4	5
4					1	2	3	4
5						1	2	3
6							1	2
7								1

*Connecting Intervals Algorithm* (Schmidt & JS, CPM 2004):

While reading  $S_2$ , mark in  $S_1$  the observed characters and track maximal intervals of marked characters.

$S_1$ 

<u>3</u>	1	2	<u>3</u>	1	<u>5</u>	2	6
0	1	2	3	4	5	6	7

$S_2$ 

4	3	5	5	5	1	4	2	2
	$i$	$j$						

## An elementary algorithm for two sequences

Preprocessing: compute two tables for  $S_1 = (3, 1, 2, 3, 1, 5, 2, 6)$ :

$POS[1] = 2, 5$   
 $POS[2] = 3, 7$   
 $POS[3] = 1, 4$   
 $POS[4] = \text{empty}$   
 $POS[5] = 6$   
 $POS[6] = 8$

$NUM(i, j)$ :

$i \setminus j$	0	1	2	3	4	5	6	7
0	1	2	3	3	3	4	4	5
1		1	2	3	3	4	4	5
2			1	2	3	4	4	5
3				1	2	3	4	5
4					1	2	3	4
5						1	2	3
6							1	2
7								1

*Connecting Intervals Algorithm* (Schmidt & JS, CPM 2004):

While reading  $S_2$ , mark in  $S_1$  the observed characters and track maximal intervals of marked characters.

$S_1$ 

<u>3</u>	1	2	<u>3</u>	1	<u>5</u>	2	6
0	1	2	3	4	5	6	7

$S_2$ 

4	3	5	5	5	1	4	2	2
	$i$		$j$					

## An elementary algorithm for two sequences

Preprocessing: compute two tables for  $S_1 = (3, 1, 2, 3, 1, 5, 2, 6)$ :

$POS[1] = 2, 5$   
 $POS[2] = 3, 7$   
 $POS[3] = 1, 4$   
 $POS[4] = \text{empty}$   
 $POS[5] = 6$   
 $POS[6] = 8$

$NUM(i, j)$ :

$i \setminus j$	0	1	2	3	4	5	6	7
0	1	2	3	3	3	4	4	5
1		1	2	3	3	4	4	5
2			1	2	3	4	4	5
3				1	2	3	4	5
4					1	2	3	4
5						1	2	3
6							1	2
7								1

*Connecting Intervals Algorithm* (Schmidt & JS, CPM 2004):

While reading  $S_2$ , mark in  $S_1$  the observed characters and track maximal intervals of marked characters.

$S_1$ 

<u>3</u>	1	2	<u>3</u>	1	<u>5</u>	2	6
0	1	2	3	4	5	6	7

$S_2$ 

4	3	5	5	5	1	4	2	2
	$i$			$j$				

## An elementary algorithm for two sequences

Preprocessing: compute two tables for  $S_1 = (3, 1, 2, 3, 1, 5, 2, 6)$ :

$POS[1] = 2, 5$   
 $POS[2] = 3, 7$   
 $POS[3] = 1, 4$   
 $POS[4] = \text{empty}$   
 $POS[5] = 6$   
 $POS[6] = 8$

$NUM(i, j)$ :

$i \setminus j$	0	1	2	3	4	5	6	7
0	1	2	3	3	3	4	4	5
1		1	2	3	3	4	4	5
2			1	2	3	4	4	5
3				1	2	3	4	5
4					1	2	3	4
5						1	2	3
6							1	2
7								1

*Connecting Intervals Algorithm* (Schmidt & JS, CPM 2004):

While reading  $S_2$ , mark in  $S_1$  the observed characters and track maximal intervals of marked characters.

$S_1$ 

<u>3</u>	<u>1</u>	2	<u>3</u>	<u>1</u>	<u>5</u>	2	6
0	1	2	3	4	5	6	7

$S_2$ 

4	3	5	5	5	1	4	2	2
	$i$				$j$			

## An elementary algorithm for two sequences

Preprocessing: compute two tables for  $S_1 = (3, 1, 2, 3, 1, 5, 2, 6)$ :

$POS[1] = 2, 5$   
 $POS[2] = 3, 7$   
 $POS[3] = 1, 4$   
 $POS[4] = \text{empty}$   
 $POS[5] = 6$   
 $POS[6] = 8$

$NUM(i, j)$ :

$i \setminus j$	0	1	2	3	4	5	6	7
0	1	2	3	3	3	4	4	5
1		1	2	3	3	4	4	5
2			1	2	3	4	4	5
3				1	2	3	4	5
4					1	2	3	4
5						1	2	3
6							1	2
7								1

*Connecting Intervals Algorithm* (Schmidt & JS, CPM 2004):

While reading  $S_2$ , mark in  $S_1$  the observed characters and track maximal intervals of marked characters.

$S_1$ 

<u>3</u>	<u>1</u>	2	<u>3</u>	<u>1</u>	<u>5</u>	2	6
0	1	2	3	4	5	6	7

$S_2$ 

4	3	5	5	5	1	4	2	2
	$i$				$j$			

Analysis:  $O(n^2)$  time and space.

## More algorithms

---

### Space reduction:

- A different algorithm based on work by Didier (WABI, 2003) finds all common intervals of two sequences in  $\mathcal{O}(n^2)$  time and  $\mathcal{O}(n)$  space.

### More than two sequences:

- Find all common intervals in  $k$  sequences in  $\mathcal{O}(kn^2)$  time and space.
- Find all common intervals that appear in at least  $k'$  out of  $k$  given sequences in  $\mathcal{O}(k(1 + k - k')n^2)$  time and  $\mathcal{O}(kn^2)$  space.



## Experimental results. Data source: COG

---

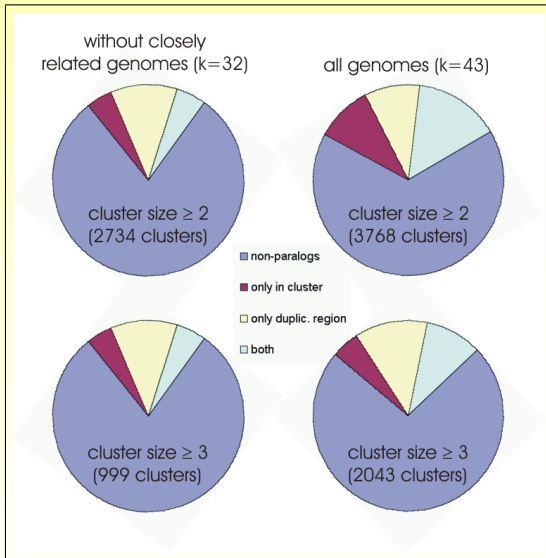
Aquifex aeolicus complete genome - 0..1551335

1529 proteins

Location	Strand	Length	PID	Synonym	Code	COG	Product
1..2100	+	699	15605613	fusA	J	COG0480	elongation factor EF-G
17..3334	+	405	15605614	tufA1	J	COG0050	elongation factor EF-Tu
46..3660	+	104	15605615	rpsJ	J	COG0051	ribosomal protein S10
3665..4390	+	241	15605616	rplC	J	COG0087	ribosomal protein L03
4387..4986	+	199	15605617	rplD	J	COG0088	ribosomal protein L04
4990..5301	+	103	15605618	rplW	J	COG0089	ribosomal protein L23
5313..6227	+	304	15605619	rplB	J	COG0090	ribosomal protein L02
6340..6900	+	186	15605620	rpsS	J	COG0185	ribosomal protein S19
7018..7314	+	98	15605621	rplV	J	COG0091	ribosomal protein L22

480 → 50 → 51 → 87 → 88 → 89

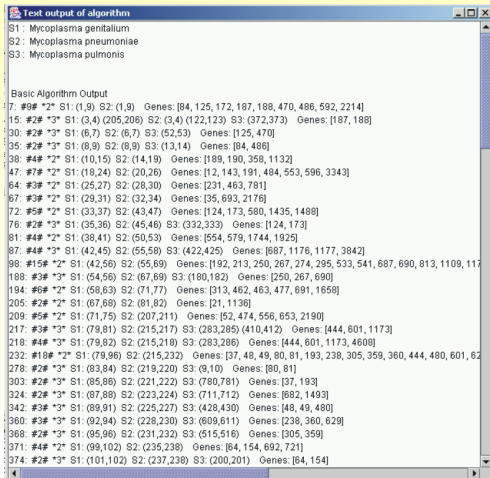
## Experimental results. Application to 43 bacterial genomes



# Experimental results. Graphical inspection of gene clusters:

**GECKO** ([bibiserv.techfak.uni-bielefeld.de/gecko](http://bibiserv.techfak.uni-bielefeld.de/gecko))

---



```
Text output of algorithm
S1: Mycoplasma genitalium
S2: Mycoplasma pneumoniae
S3: Mycoplasma pulmonis

Basic Algorithm Output
7: #9# *2* S1: (1,9) S2: (1,9) Genes: [84, 125, 172, 187, 188, 470, 486, 592, 2214]
15: #2# *3* S1: (3,4) (205,206) S2: (3,4) (122,123) S3: (372,373) Genes: [187, 188]
30: #2# *3* S1: (6,7) S2: (6,7) S3: (52,53) Genes: [125, 470]
35: #2# *3* S1: (8,9) S2: (8,9) S3: (13,14) Genes: [84, 486]
38: #4# *2* S1: (10,15) S2: (14,19) Genes: [189, 190, 358, 1132]
47: #7# *2* S1: (18,24) S2: (20,26) Genes: [12, 143, 191, 484, 553, 596, 3343]
64: #3# *2* S1: (25,27) S2: (28,30) Genes: [231, 463, 781]
67: #3# *2* S1: (29,31) S2: (32,34) Genes: [35, 693, 2176]
72: #5# *2* S1: (33,37) S2: (43,47) Genes: [124, 173, 580, 1435, 1488]
76: #2# *3* S1: (35,36) S2: (45,46) S3: (332,333) Genes: [124, 173]
81: #4# *2* S1: (38,41) S2: (50,53) Genes: [554, 579, 1744, 1925]
87: #4# *3* S1: (42,45) S2: (55,58) S3: (422,425) Genes: [687, 1176, 1177, 3842]
98: #15# *2* S1: (42,56) S2: (55,69) Genes: [192, 213, 250, 267, 274, 295, 533, 541, 687, 690, 813, 1109, 1177]
188: #3# *3* S1: (54,56) S2: (67,69) S3: (180,182) Genes: [250, 267, 690]
194: #6# *2* S1: (58,63) S2: (71,77) Genes: [313, 462, 463, 477, 691, 1658]
205: #2# *2* S1: (67,68) S2: (81,82) Genes: [21, 1136]
209: #5# *2* S1: (71,75) S2: (207,211) Genes: [52, 474, 556, 653, 2190]
217: #3# *3* S1: (79,81) S2: (215,217) S3: (283,285) (410,412) Genes: [444, 601, 1173]
218: #4# *3* S1: (79,82) S2: (215,218) S3: (283,286) Genes: [444, 601, 1173, 4608]
232: #18# *2* S1: (79,96) S2: (215,232) Genes: [37, 48, 49, 80, 81, 193, 238, 305, 359, 360, 444, 480, 601, 629]
278: #2# *3* S1: (83,84) S2: (219,220) S3: (9,10) Genes: [80, 81]
303: #2# *3* S1: (85,86) S2: (221,222) S3: (780,781) Genes: [37, 193]
324: #2# *3* S1: (87,88) S2: (223,224) S3: (711,712) Genes: [682, 1493]
342: #3# *3* S1: (89,91) S2: (225,227) S3: (428,430) Genes: [48, 49, 480]
360: #3# *3* S1: (92,94) S2: (228,230) S3: (609,611) Genes: [238, 360, 629]
368: #2# *3* S1: (95,96) S2: (231,232) S3: (515,516) Genes: [305, 359]
371: #4# *2* S1: (99,102) S2: (235,238) Genes: [64, 154, 692, 721]
374: #2# *3* S1: (101,102) S2: (237,238) S3: (200,201) Genes: [64, 154]
```

# Experimental results. Graphical inspection of gene clusters:

## GECKO (bibiserv.techfak.uni-bielefeld.de/gecko)

The screenshot displays the GeneCluster V1.0\_pre1 software interface. The main window shows a list of clusters with columns for ID, #Genes, #Seq, #SubCl, and Contained Genes. The clusters are sorted by the number of genes. The graphical representation on the right shows the gene clusters as colored arrows pointing to the right, with the gene IDs (735, 803, 1108, 1121, 1108, 886) labeled above each arrow. The clusters are grouped by species: Thermotoga maritima, Streptococcus pneumoniae, Lactococcus lactis subsp. lacti, Listeria innocua Clip11262, Synechocystis sp. PCC 6803, and Pseudomonas aeruginosa PA01.

**Text output of algorithm**

S1: Mycoplasma genitalium  
S2: Mycoplasma pneumoniae  
S3: Mycoplasma pulmonis

**Basic Algorithm Output**

7: #9# \*2\* S1: (1,9) S2: (1,9)  
15: #2# \*3\* S1: (3,4) S2: (205,201)  
30: #2# \*3\* S1: (6,7) S2: (6,7)  
35: #2# \*3\* S1: (8,9) S2: (8,9)  
38: #4# \*2\* S1: (10,15) S2: (10,15)  
47: #7# \*2\* S1: (18,24) S2: (18,24)  
64: #3# \*2\* S1: (25,27) S2: (25,27)  
67: #3# \*2\* S1: (29,31) S2: (29,31)  
72: #5# \*2\* S1: (33,37) S2: (33,37)  
76: #2# \*3\* S1: (35,36) S2: (35,36)  
81: #4# \*2\* S1: (38,41) S2: (38,41)  
87: #4# \*3\* S1: (42,45) S2: (42,45)  
98: #15# \*2\* S1: (42,56) S2: (42,56)  
188: #3# \*3\* S1: (54,56) S2: (54,56)  
194: #6# \*2\* S1: (58,63) S2: (58,63)  
205: #2# \*2\* S1: (67,68) S2: (67,68)  
209: #5# \*2\* S1: (71,75) S2: (71,75)  
217: #3# \*3\* S1: (79,81) S2: (79,81)  
218: #4# \*3\* S1: (79,82) S2: (79,82)  
232: #18# \*2\* S1: (79,96) S2: (79,96)  
278: #2# \*3\* S1: (83,84) S2: (83,84)  
303: #2# \*3\* S1: (85,86) S2: (85,86)  
324: #2# \*3\* S1: (87,88) S2: (87,88)  
342: #3# \*3\* S1: (89,91) S2: (89,91)  
360: #3# \*3\* S1: (92,94) S2: (228,230) S3: (609,611) Genes: [238, 360, 629]  
368: #2# \*3\* S1: (95,96) S2: (231,232) S3: (515,516) Genes: [305, 359]  
371: #4# \*2\* S1: (99,102) S2: (235,238) Genes: [64, 154, 692, 721]  
374: #2# \*3\* S1: (101,102) S2: (237,238) S3: (200,201) Genes: [64, 154]

ID	#Genes	#Seq	#SubCl	Contained Genes
274	5	3	0	[195, 532, 779, 1358, 2740]
284	5	3	0	[134, 135, 147, 512, 547]
292	5	4	2	[48, 49, 50, 51, 480]
293	5	4	2	[444, 601, 747, 1124, 1173]
312	5	6	2	[735, 803, 1108, 1121, 1846]
14	6	3	0	[34, 46, 138, 150, 151, 299]
279	6	3	0	[206, 325, 762, 849, 1799, 2302]
7	6	4	0	[81, 81, 777, 744, 750, 690]

ID	#Genes	#Seq	Contained Genes
1	4	3	[735, 803, 1108, 1121]
288	4	3	[803, 1108, 1121, 1846]

ID	#Genes	#Seq	Contained Genes
1	4	3	[735, 803, 1108, 1121]
288	4	3	[803, 1108, 1121, 1846]

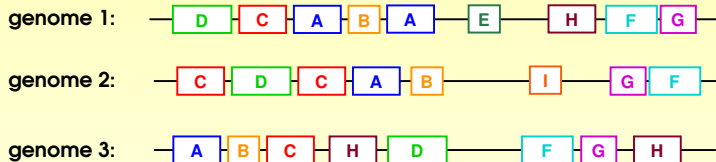
Thermotoga maritima 735 803 1121 1108  
Streptococcus pneumoniae 803 1108 1121 886  
Lactococcus lactis subsp. lacti 1108 1121 803 886  
Listeria innocua Clip11262 803 1108 1121 886  
Synechocystis sp. PCC 6803 735 803 1121 1108  
Pseudomonas aeruginosa PA01 803 735 1121 1108

Thermotoga maritima 735 803 1121 1108  
Synechocystis sp. PCC 6803 735 803 1121 1108  
Pseudomonas aeruginosa PA01 803 735 1121 1108

Streptococcus pneumoniae 803 1108 1121 886  
Lactococcus lactis subsp. lacti 1108 1121 803 886  
Listeria innocua Clip11262 803 1108 1121 886

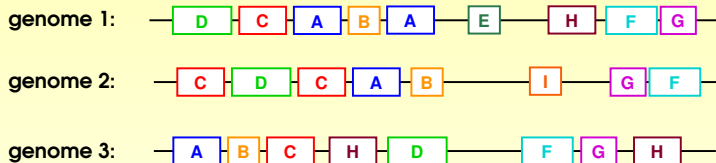
## Approximate Gene Clusters

---



## Approximate Gene Clusters

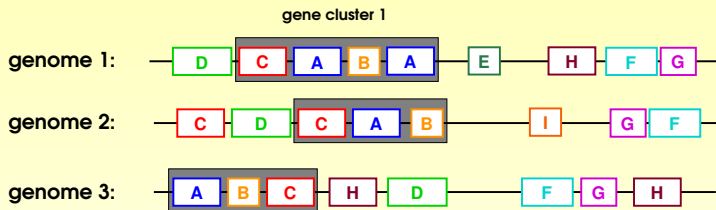
---



- (perfect) gene cluster:
  - set of genes occurring *en bloc* in multiple genomes
  - gene order within blocks not considered
  - multiple occurrences of genes possible

# Approximate Gene Clusters

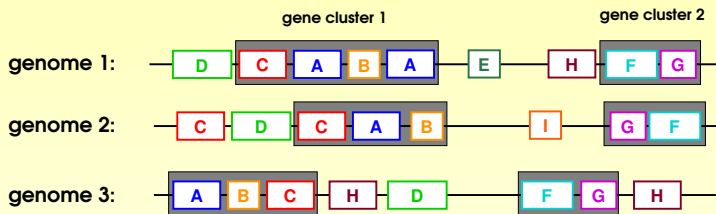
---



- (perfect) gene cluster:
  - set of genes occurring *en bloc* in multiple genomes
  - gene order within blocks not considered
  - multiple occurrences of genes possible

# Approximate Gene Clusters

---

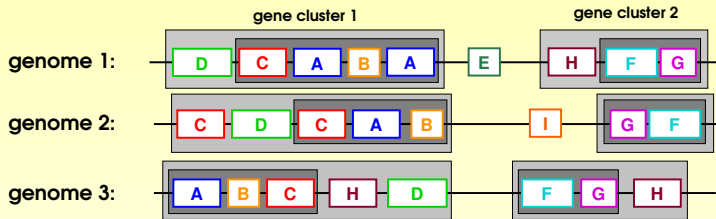


- (perfect) gene cluster:
  - set of genes occurring *en bloc* in multiple genomes
  - gene order within blocks not considered
  - multiple occurrences of genes possible



# Approximate Gene Clusters

---



- (perfect) gene cluster:
  - set of genes occurring *en bloc* in multiple genomes
  - gene order within blocks not considered
  - multiple occurrences of genes possible
- approximate gene cluster:
  - additional and missing genes

## Models for Approximate Gene Clusters

---

- ***r*-window model** (Friedman and Hughes, 2001)
  - fixed block size
  - pairwise comparison in polynomial time
- **max-gap model** (Bergeron, Corteel and Raffinot, 2002)
  - upper bound for insertion length
  - pairwise comparison in polynomial time
  - exponential in number of compared sequences
- **approximate gene clusters** (Rahmann and Klau, 2006)
  - very general cluster model, subsumes most other models
  - ILP approach
- **median gene clusters** (Böcker, Jahn, Mixtacki and JS, 2008)

# Median of Character Sets

---

Sequence 1: ..... 1 2 3 4 5 6 7 .....  
Sequence 2: ... 1 8 3 4 5 2 1 6 7 9 .....  
Sequence 3: ..... 8 3 7 5 1 .....  
Sequence 4: ... 3 1 4 3 .....  
Sequence 5: ..... 2 4 5 9 2 7 3 ...

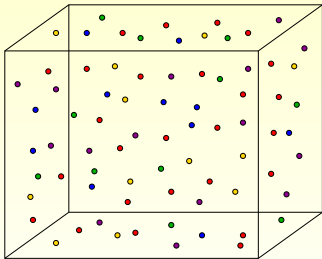
$\Sigma$ :	1	2	3	4	5	6	7	8	9
$CS(S_1(i_1, j_1))$	1	1	1	1	1	1	1	0	0
$CS(S_2(i_2, j_2))$	1	1	1	1	1	1	1	1	1
$CS(S_3(i_3, j_3))$	1	0	1	0	1	0	1	1	0
$CS(S_4(i_4, j_4))$	1	0	1	1	0	0	0	0	0
$CS(S_5(i_5, j_5))$	0	1	1	1	1	0	1	0	1

# Median of Character Sets

Sequence 1: ..... 1 2 3 4 5 6 7 .....  
Sequence 2: ... 1 8 3 4 5 2 1 6 7 9 .....  
Sequence 3: ..... 8 3 7 5 1 .....  
Sequence 4: ... 3 1 4 3 .....  
Sequence 5: ..... 2 4 5 9 2 7 3 ...

$\Sigma :$	1	2	3	4	5	6	7	8	9
$CS(S_1(i_1, j_1))$	1	1	1	1	1	1	1	0	0
$CS(S_2(i_2, j_2))$	1	1	1	1	1	1	1	1	1
$CS(S_3(i_3, j_3))$	1	0	1	0	1	0	1	1	0
$CS(S_4(i_4, j_4))$	1	0	1	1	0	0	0	0	0
$CS(S_5(i_5, j_5))$	0	1	1	1	1	0	1	0	1

Space of Character Sets:  $\{0, 1\}^{|\Sigma|}$

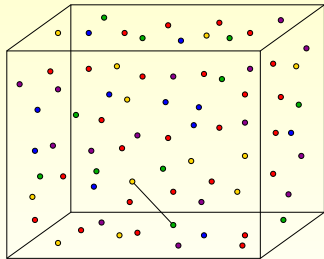


# Median of Character Sets

Sequence 1: ..... 1 2 3 4 5 6 7 .....  
 Sequence 2: ... 1 8 3 4 5 2 1 6 7 9 .....  
 Sequence 3: ..... 8 3 7 5 1 .....  
 Sequence 4: ... 3 1 4 3 .....  
 Sequence 5: ..... 2 4 5 9 2 7 3 ...

$\Sigma :$	1	2	3	4	5	6	7	8	9
$CS(S_1(i_1, j_1))$	1	1	1	1	1	1	1	0	0
$CS(S_2(i_2, j_2))$	1	1	1	1	1	1	1	1	1
$CS(S_3(i_3, j_3))$	1	0	1	0	1	0	1	1	0
$CS(S_4(i_4, j_4))$	1	0	1	1	0	0	0	0	0
$CS(S_5(i_5, j_5))$	0	1	1	1	1	0	1	0	1

Space of Character Sets:  $\{0, 1\}^{|\Sigma|}$

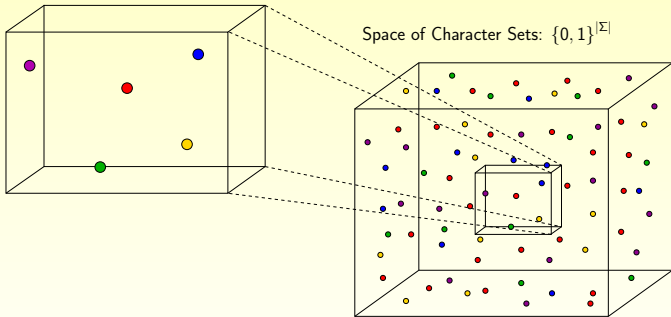


$$D_{sym}(C_1, C_2) = |C_1 \setminus C_2| + |C_2 \setminus C_1|$$

# Median of Character Sets

Sequence 1: ..... 1 2 3 4 5 6 7 .....  
 Sequence 2: ... 1 8 3 4 5 2 1 6 7 9 .....  
 Sequence 3: ..... 8 3 7 5 1 .....  
 Sequence 4: ... 3 1 4 3 .....  
 Sequence 5: ..... 2 4 5 9 2 7 3 ...

$\Sigma :$	1	2	3	4	5	6	7	8	9
$CS(S_1(i_1, j_1))$	1	1	1	1	1	1	1	0	0
$CS(S_2(i_2, j_2))$	1	1	1	1	1	1	1	1	1
$CS(S_3(i_3, j_3))$	1	0	1	0	1	0	1	1	0
$CS(S_4(i_4, j_4))$	1	0	1	1	0	0	0	0	0
$CS(S_5(i_5, j_5))$	0	1	1	1	1	0	1	0	1

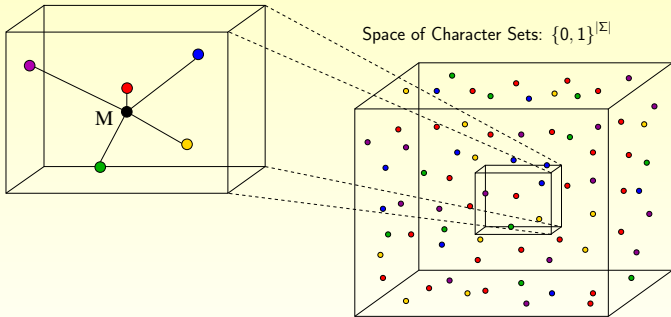


$$D_{sym}(C_1, C_2) = |C_1 \setminus C_2| + |C_2 \setminus C_1|$$

# Median of Character Sets

Sequence 1: ..... 1 2 3 4 5 6 7 .....  
 Sequence 2: ... 1 8 3 4 5 2 1 6 7 9 .....  
 Sequence 3: ..... 8 3 7 5 1 .....  
 Sequence 4: ... 3 1 4 3 .....  
 Sequence 5: ..... 2 4 5 9 2 7 3 ...

$\Sigma :$	1	2	3	4	5	6	7	8	9
$CS(S_1(i_1, j_1))$	1	1	1	1	1	1	1	0	0
$CS(S_2(i_2, j_2))$	1	1	1	1	1	1	1	1	1
$CS(S_3(i_3, j_3))$	1	0	1	0	1	0	1	1	0
$CS(S_4(i_4, j_4))$	1	0	1	1	0	0	0	0	0
$CS(S_5(i_5, j_5))$	0	1	1	1	1	0	1	0	1

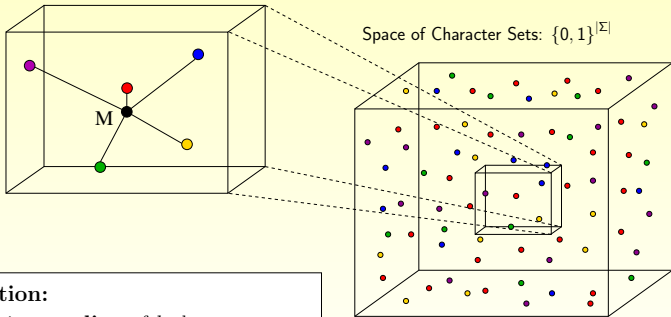


$$D_{sym}(C_1, C_2) = |C_1 \setminus C_2| + |C_2 \setminus C_1|$$

# Median of Character Sets

Sequence 1: ..... 1 2 3 4 5 6 7 .....  
 Sequence 2: ... 1 8 3 4 5 2 1 6 7 9 .....  
 Sequence 3: ..... 8 3 7 5 1 .....  
 Sequence 4: ... 3 1 4 3 .....  
 Sequence 5: ..... 2 4 5 9 2 7 3 ...

$\Sigma :$	1	2	3	4	5	6	7	8	9
$CS(S_1(i_1, j_1))$	1	1	1	1	1	1	1	0	0
$CS(S_2(i_2, j_2))$	1	1	1	1	1	1	1	1	1
$CS(S_3(i_3, j_3))$	1	0	1	0	1	0	1	1	0
$CS(S_4(i_4, j_4))$	1	0	1	1	0	0	0	0	0
$CS(S_5(i_5, j_5))$	0	1	1	1	1	0	1	0	1



## Definition:

$M \subseteq \Sigma$  is a **median** of  $k$  character sets  $C_1, \dots, C_k \subseteq \Sigma$  iff for all  $C \subseteq \Sigma$ :

$$\sum_{\ell=1}^k D_{sym}(M, C_\ell) \leq \sum_{\ell=1}^k D_{sym}(C, C_\ell)$$

$$D_{sym}(C_1, C_2) = |C_1 \setminus C_2| + |C_2 \setminus C_1|$$



## Problem Statement

---

### Given:

- sequences  $S_1, \dots, S_k$  over alphabet  $\Sigma$
- $s$  (minimum cluster size)
- $\delta$  (distance threshold)

### Wanted: all $M \subseteq \Sigma$ with

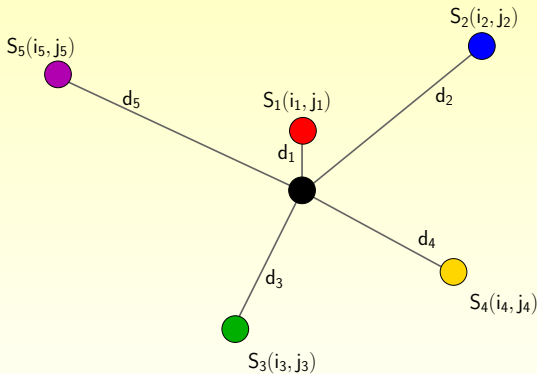
- $M$  is a median for some  $S_1[i_1, j_1], \dots, S_k[i_k, j_k]$
- $\sum_{\ell=1}^k D(M, \text{CS}(S[i_\ell, j_\ell])) \leq \delta$
- $|M| \geq s$

Such a set  $M$  is called a **median gene cluster** of  $S_1, \dots, S_k$ .

## Main Idea: Search Space Reduction

---

- search space: all  $\mathcal{O}(n^{2k})$  combinations of substrings of  $S_1, \dots, S_k$
- cluster filter approach:

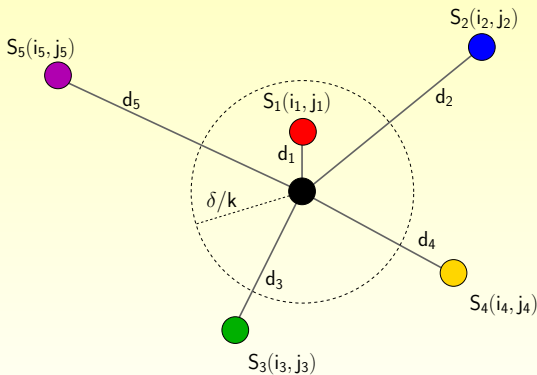


- median distance threshold:  $\sum_{\ell=1}^k d_{\ell} \leq \delta$

## Main Idea: Search Space Reduction

---

- search space: all  $\mathcal{O}(n^{2k})$  combinations of substrings of  $S_1, \dots, S_k$
- cluster filter approach:

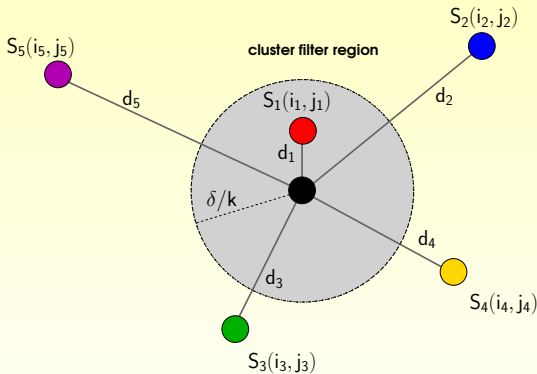


- median distance threshold:  $\sum_{\ell=1}^k d_{\ell} \leq \delta$

## Main Idea: Search Space Reduction

---

- search space: all  $\mathcal{O}(n^{2k})$  combinations of substrings of  $S_1, \dots, S_k$
- cluster filter approach:

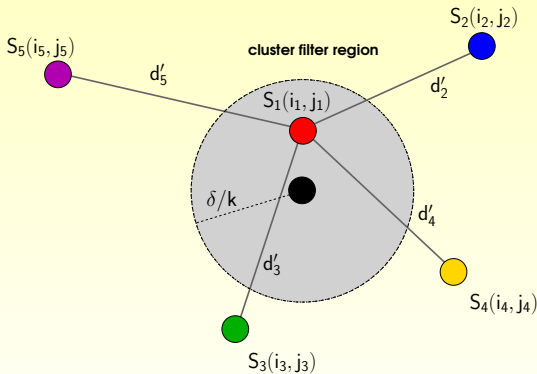


- median distance threshold:  $\sum_{\ell=1}^k d_{\ell} \leq \delta$

## Main Idea: Search Space Reduction

---

- search space: all  $\mathcal{O}(n^{2k})$  combinations of substrings of  $S_1, \dots, S_k$
- cluster filter approach:



- median distance threshold:  $\sum_{\ell=1}^k d_{\ell} \leq \delta$
- cluster filter distance threshold:  $\sum_{\ell=1}^k d'_{\ell} \leq 2 \frac{k-1}{k} \delta$

### 3-Step Approach to Median Gene Cluster Computation

- **Step 1:** compute the set of all cluster filters  $C$  for  $S_1, \dots, S_k$
- **Step 2:** compute for each cluster filter  $C$  all combinations with substrings from the other sequences for which:

$$\sum_{\ell=1}^k D(C, \mathcal{CS}(S_\ell[i_\ell, j_\ell])) \leq 2 \frac{k-1}{k} \delta$$

- **Step 3:**
  - computation of median(s) of each  $k$ -tuple (from Step 2)
  - comparison of median distance with distance threshold

## Step 1: Computation of Cluster Filters

---

- naive approach: testing all  $\mathcal{O}(k^2n^4)$  combinations of substrings
- our approach: extension of the *Connecting Intervals Algorithm* (Schmidt and JS, 2004)

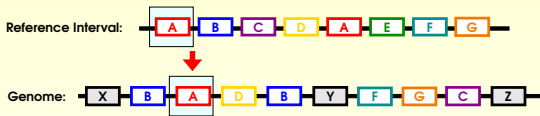
Reference Interval: 

Genome: 

## Step 1: Computation of Cluster Filters

---

- naive approach: testing all  $\mathcal{O}(k^2n^4)$  combinations of substrings
- our approach: extension of the *Connecting Intervals Algorithm* (Schmidt and JS, 2004)

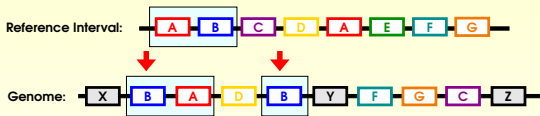




## Step 1: Computation of Cluster Filters

---

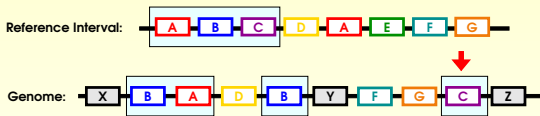
- naive approach: testing all  $\mathcal{O}(k^2n^4)$  combinations of substrings
- our approach: extension of the *Connecting Intervals Algorithm* (Schmidt and JS, 2004)



## Step 1: Computation of Cluster Filters

---

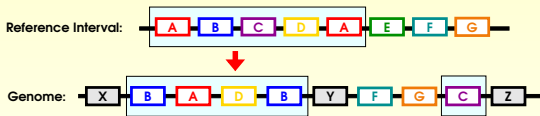
- naive approach: testing all  $\mathcal{O}(k^2n^4)$  combinations of substrings
- our approach: extension of the *Connecting Intervals Algorithm* (Schmidt and JS, 2004)



## Step 1: Computation of Cluster Filters

---

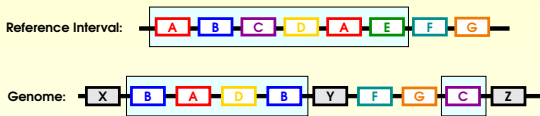
- naive approach: testing all  $\mathcal{O}(k^2n^4)$  combinations of substrings
- our approach: extension of the *Connecting Intervals Algorithm* (Schmidt and JS, 2004)



## Step 1: Computation of Cluster Filters

---

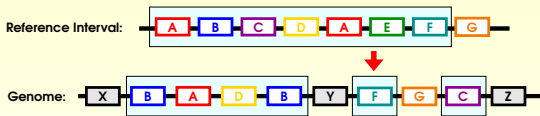
- naive approach: testing all  $\mathcal{O}(k^2n^4)$  combinations of substrings
- our approach: extension of the *Connecting Intervals Algorithm* (Schmidt and JS, 2004)



## Step 1: Computation of Cluster Filters

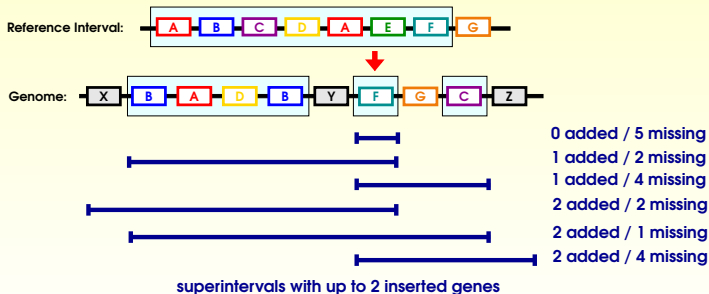
---

- naive approach: testing all  $\mathcal{O}(k^2n^4)$  combinations of substrings
- our approach: extension of the *Connecting Intervals Algorithm* (Schmidt and JS, 2004)



## Step 1: Computation of Cluster Filters

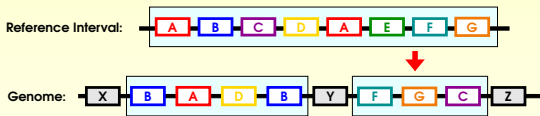
- naive approach: testing all  $\mathcal{O}(k^2n^4)$  combinations of substrings
- our approach: extension of the *Connecting Intervals Algorithm* (Schmidt and JS, 2004)



## Step 1: Computation of Cluster Filters

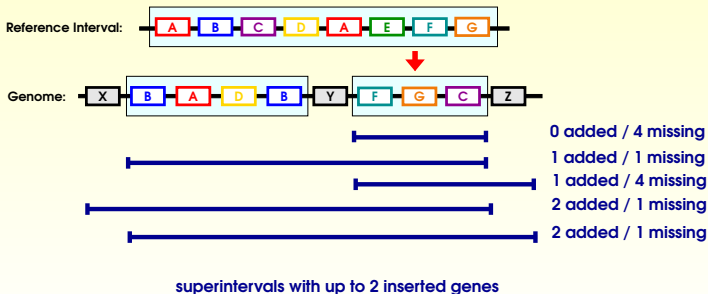
---

- naive approach: testing all  $\mathcal{O}(k^2n^4)$  combinations of substrings
- our approach: extension of the *Connecting Intervals Algorithm* (Schmidt and JS, 2004)



## Step 1: Computation of Cluster Filters

- naive approach: testing all  $\mathcal{O}(k^2n^4)$  combinations of substrings
- our approach: extension of the *Connecting Intervals Algorithm* (Schmidt and JS, 2004)

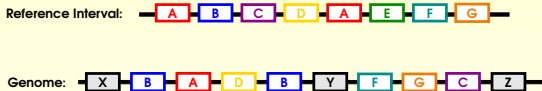




## Step 1: Computation of Cluster Filters

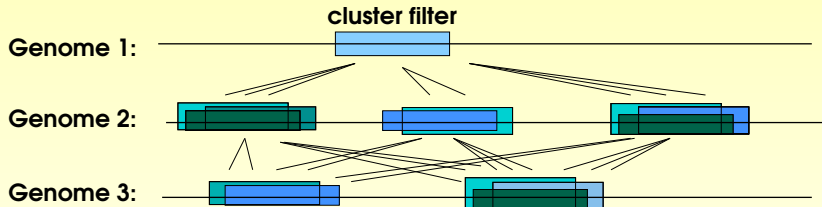
---

- naive approach: testing all  $\mathcal{O}(k^2n^4)$  combinations of substrings
- our approach: extension of the *Connecting Intervals Algorithm* (Schmidt and JS, 2004)



- runtime:  $\mathcal{O}(k^2n^2(1 + \delta^2))$ , space:  $\mathcal{O}(kn^2)$

## Step 2: Computation of k-Tuples Containing a Cluster Filter



**Idea:** Build for each cluster filter  $C$  all combinations with substrings of the other sequences for which  $C$  is cluster filter

- possible combinatorial explosion  $\mathcal{O}(n^k)$  (unlikely in practice)
- up to  $\mathcal{O}(\delta^{2k})$  variants of a  $k$ -tuple

### Step 3: Computation of Median of Each k-Tuple

---

- majority vote:  $\mathcal{O}(k|\Sigma|)$  time and  $\mathcal{O}(k|\Sigma|)$  space
- compare total distance with threshold  $\delta$

	A	B	C	D	E	F	G	H	I	J
Genome 1	1	1	1	1	0	1	0	1	1	0
Genome 2	1	1	0	1	0	1	1	1	1	1
Genome 3	1	1	1	1	1	1	0	1	0	1
Genome 4	1	1	1	1	1	1	1	1	1	0
Genome 5	1	0	1	1	0	1	0	1	1	1
	5	4	4	5	2	5	2	5	4	3

Median = { A B C D F H I J }

### Step 3: Computation of Median of Each k-Tuple

- majority vote:  $\mathcal{O}(k|\Sigma|)$  time and  $\mathcal{O}(k|\Sigma|)$  space
- compare total distance with threshold  $\delta$

	A	B	C	D	E	F	G	H	I	J
Genome 1	1	1	1	1	0	1	0	1	1	0
Genome 2	1	1	0	1	0	1	1	1	1	1
Genome 3	1	1	1	1	1	1	0	1	0	1
Genome 4	1	1	1	1	1	1	1	1	1	0
Genome 5	1	0	1	1	0	1	0	1	1	1
	5	4	4	5	2	5	2	5	4	3

Median = { A B C D F H I J }

## Algorithm Summary

---

- **Step 1:** cluster filter computation  $\mathcal{O}(k^2 n^2 (1 + \delta^2))$
- **Step 2:** combination of each cluster filter  $C$  with substrings from other sequences with

$$\sum_{\ell=1}^k D(C, \mathcal{CS}(S_\ell[i_\ell, j_\ell])) \leq 2 \frac{k-1}{k} \delta$$

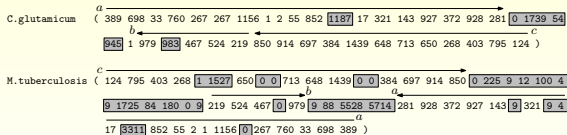
possible combinatorial explosion:  $\mathcal{O}(n^k)$

- **Step 3:**
  - computation of median of each such  $k$ -tuple  $\mathcal{O}(k|\Sigma|)$
  - comparison: median distance vs. distance threshold  $\delta$

## Experimental Results

---

- comparison with ILP approach (Rahmann and Klau 2006)
- dataset: genomes of *C. glutamicum* and *M. tuberculosis*
- detection of gene cluster containing 51 genes
- runtime of ILP program: > 1h
- runtime of our approach: 17 sec. (on slower processor)



## Experimental Results

---

- application to gene cluster detection in multiple genomes
- five  $\gamma$ -proteobacteria:
  - *Buchnera aphidicola* APS (NC\_002528)
  - *Escherichia coli* K12 (NC\_000913)
  - *Haemophilus influenzae* Rd (NC\_000907)
  - *Pasteurella multocida* Pm70 (NC\_002663)
  - *Xylella fastidiosa* 9a5c (NC\_002488)

	$\delta=0$	$\delta=1$	$\delta=5$	$\delta=8$	$\delta=12$
$s=10$	6	124	252	329	406
$s=15$	0	42	127	165	181
$s=25$	0	0	10	10	10

# median gene clusters

	$\delta=0$	$\delta=1$	$\delta=5$	$\delta=8$	$\delta=12$
$s=10$	6.9	7.4	20.0	103.4	1610.0
$s=15$	6.9	7.3	15.1	59.7	1148.5
$s=25$	6.9	7.1	10.5	17.6	630.0

computation time (seconds)

# Experimental Results

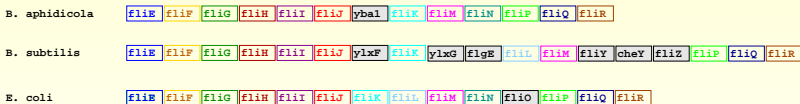
---

Sample gene clusters detected by our method:

- ATP synthesis



- flagellar biosynthesis



- cell division and cell wall biosynthesis





## Algorithmic results

---

- Find all *common intervals* of  $k$  permutations in  $\mathcal{O}(kn + |\text{output}|)$  time.
- Find all *common intervals* of  $k$  sequences in  $\mathcal{O}(kn^2)$  time.
- Find all *median gene clusters* of  $k$  sequences in  $\mathcal{O}(k^2n^2(1 + \delta^2) + n^k + k|\Sigma|)$  time

## Conclusion

---

### Points raised:

- Comparative genomics can help in functional genome annotation
- Conserved regions in genomes have a static and a dynamic aspect
- Interesting combinatorics in Bioinformatics

### Next steps:

- Statistical assessment of gene clusters
- Patterns in overlapping gene clusters
- Application to more data

## Acknowledgments

---

- Sebastian Böcker (Bielefeld/Jena)
- Steffen Heber (Raleigh)
- **Katharina Jahn** (Bielefeld)
- Hannes Luz (Berlin)
- Julia Mixtacki (Bielefeld)
- Mathieu Raffinot (Paris/Moscow)
- Sven Rahmann (Bielefeld/Dortmund)
- **Thomas Schmidt** (Bielefeld/München)

### References:

- Heber & JS, Finding all Common Intervals of  $k$  Permutations, Proc. CPM 2001.
- Heber & JS, Algorithms for Finding Gene Clusters, Proc. WABI 2001.
- Schmidt & JS, Quadratic Time Algorithms for Finding Common Intervals in Two and More Sequences, Proc. CPM 2004.
- Böcker, Jahn, Mixtacki & JS, Computation of Median Gene Clusters, Proc. RECOMB 2008.

