

Common intervals of genomes

Mathieu Raffinot
CNRS - LIAFA

Context:

- comparative genomics.
- set of genomes partially/totally annotated

➡ Informative group of genes or domains ?

Ex: COG database

Fichier Édition Affichage Historique Marque-pages Outils Aide

http://www.ncbi.nlm.nih.gov/COG/grace/shokog.cgi?KOG0383

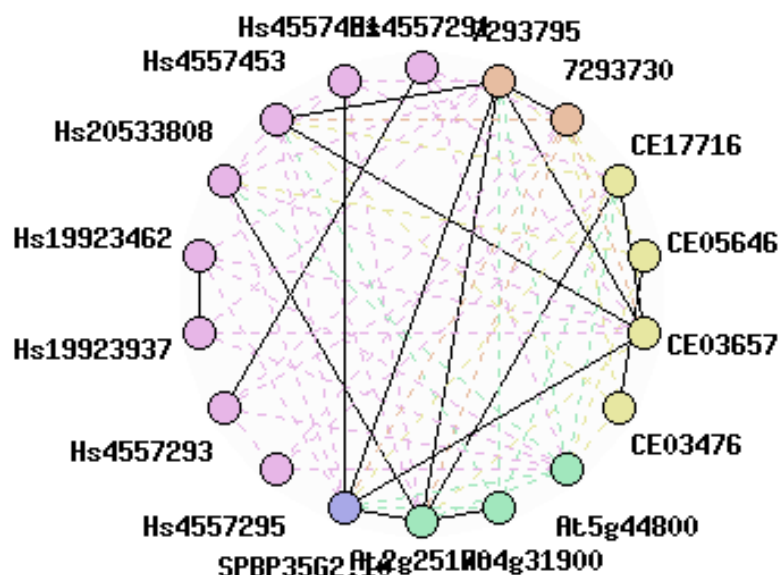
Google Yahoo! France Libération.fr Яндекс CNRS Bioinformation :: Index Les trains russes Le Monde.fr : A la une Vivre en Russie :: In... Dictionnaire anglais Mathieu Raffinot

Disable Cookies CSS Forms Images Information Miscellaneous Outline Resize Tools View Source Options

FoxLingo Page Web Texte Services TradAuto Recherche

	KOGs	TWOGs	LSEs	Kognitor			
	ACDH-P- <i>Arabidopsis thaliana</i>	R KOG0383 <i>Caenorhabditis elegans</i>	<i>Drosophila melanogaster</i>	<i>Homo sapiens</i>	<i>Sacchar. cerevisiae</i>	<i>Schizosac. pombe</i>	<i>Enceph. cuciculi</i>
	At2g25170 At4g31900 At5g44800	CE03476 CE03657 CE05646 CE17716	7293730 7293795	Hs4557291 Hs4557451 Hs4557453 Hs20533808 Hs19923462 Hs19923937 Hs4557293 Hs4557295	-	SPBP35G2.10	-

BeTs for KOG0383:



[At2g25170]	CE17716	7293795	Hs20533808	YER164w	SPAC3G6.01	ECU01g0350
[At4g31900]	CE17716	7293795	Hs4557453	YER164w	SPAC3G6.01	ECU10g1320
[At5g44800]	CE17716	7293795	Hs20533808	YER164w	SPAC3G6.01	ECU01g0350
At2g36720	[CE03476]	7293618	Hs4557453	YMR075w	SPAC16C9.05	-
At2g25170	[CE03657]	7293795	Hs4557453	YER164w	SPAC3G6.01	ECU01g0350
At3g14980	[CE05646]	7293730	Hs20533808	YMR075w	SPAC2F7.07c	-
At2g25170	[CE17716]	7293795	Hs4557453	YER164w	SPAC3G6.01	ECU01g0350
At2g25170	CE03657	[7293730]	Hs4557453	YER164w	SPAC1783.05	ECU01g0350
At2g25170	CE03657	[7293795]	Hs4557453	YER164w	SPAC3G6.01	ECU01g0350
At5g44800	CE03657	7293730	[Hs4557291]	YMR075w	SPAC16C9.05	-
At2g25170	CE17716	7293795	[Hs4557451]	YER164w	SPAC3G6.01	ECU01g0350
At2g25170	CE03657	7293795	[Hs4557453]	YER164w	SPAC3G6.01	ECU01g0350
At2g25170	CE03657	7293795	[Hs20533808]	YER164w	SPAC3G6.01	ECU01g0350
At5g44800	CE03657	7293795	[Hs19923462]	YMR075w	SPAC2F7.07c	ECU08g0060
At5g44800	CE03657	7293795	[Hs19923937]	YMR075w	SPAC2F7.07c	-
At5g44800	CE17716	7293730	[Hs4557293]	YMR075w	SPAC16C9.05	-
At5g44800	CE17716	7293730	[Hs4557295]	YMR075w	SPAC16C9.05	-
At2g25170	CE03657	7293795	Hs4557451	YER164w	[SPBP35G2.10]	ECU01g0350

Many difficulties !

Biology

What are two similar genes ? What about alternative splicing ?

When are two genes close (notion of distance) ?

What is an interesting cluster ?

basis: pressure selection -> keep genes working together close

How to model clusters ? Graphs / strings ?

How to compute those clusters ?

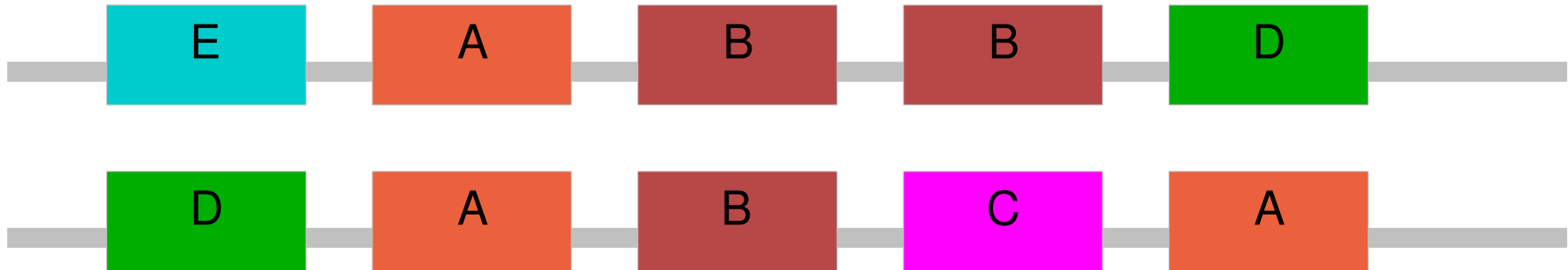
How to manage the sets of clusters and extract useful information ?

Computer science

One of the simplest model :

- genomes as strings of units
- common intervals

Simplest case in this model: 2 genomes !

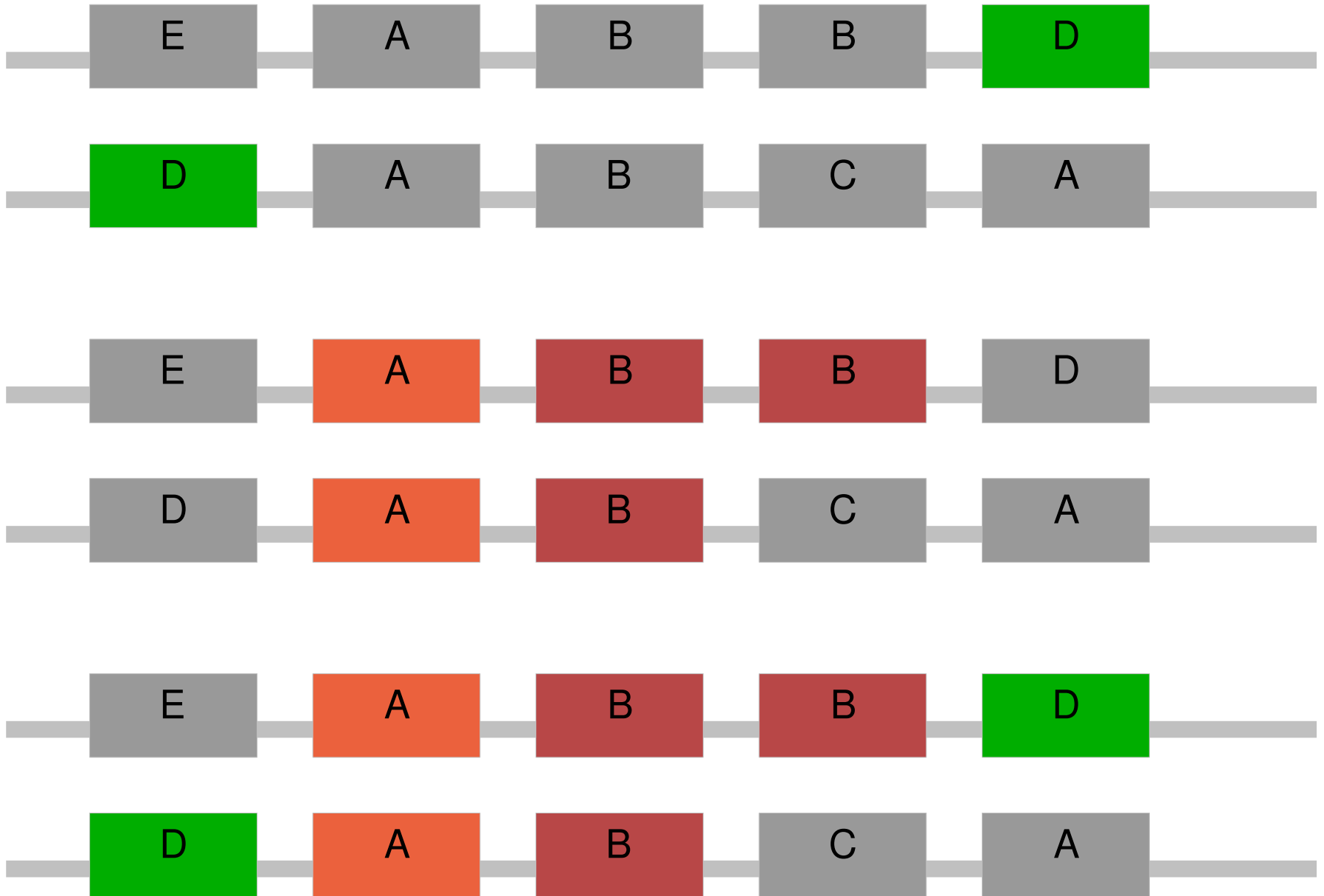


Common interval:

- one interval on each chromosome
- same set of gene in each interval
- external bounds not in the set of gene



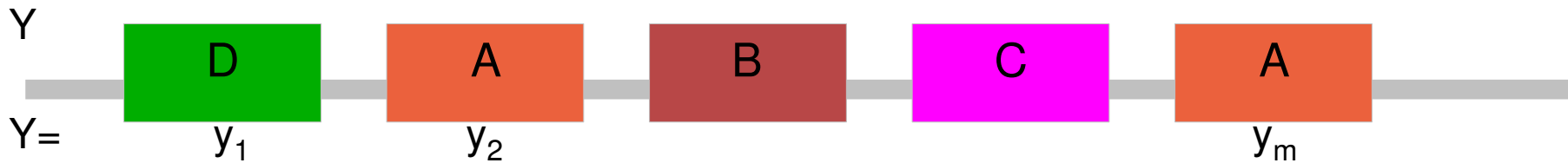




How many common intervals ?

- X first chromosome, $X = x_1 x_2 \dots x_n$
- Y second chromosome, $Y = y_1 y_2 \dots y_m$

Common alphabet Σ , $|\Sigma| \leq \max(|X|, |Y|)$



fo(Y,1) = D A B C

fo(Y,2) = A B C

fo(Y,3) = B C A

fo(Y,4) = C A

fo(Y,5) = A

D = 1 A = 2 B = 3 C = 4

A = 1 B = 2 C = 3

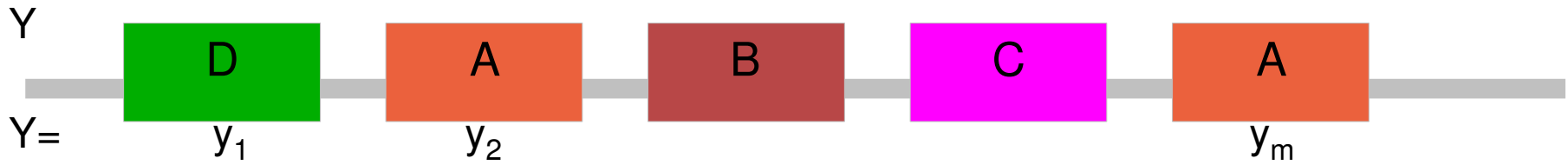
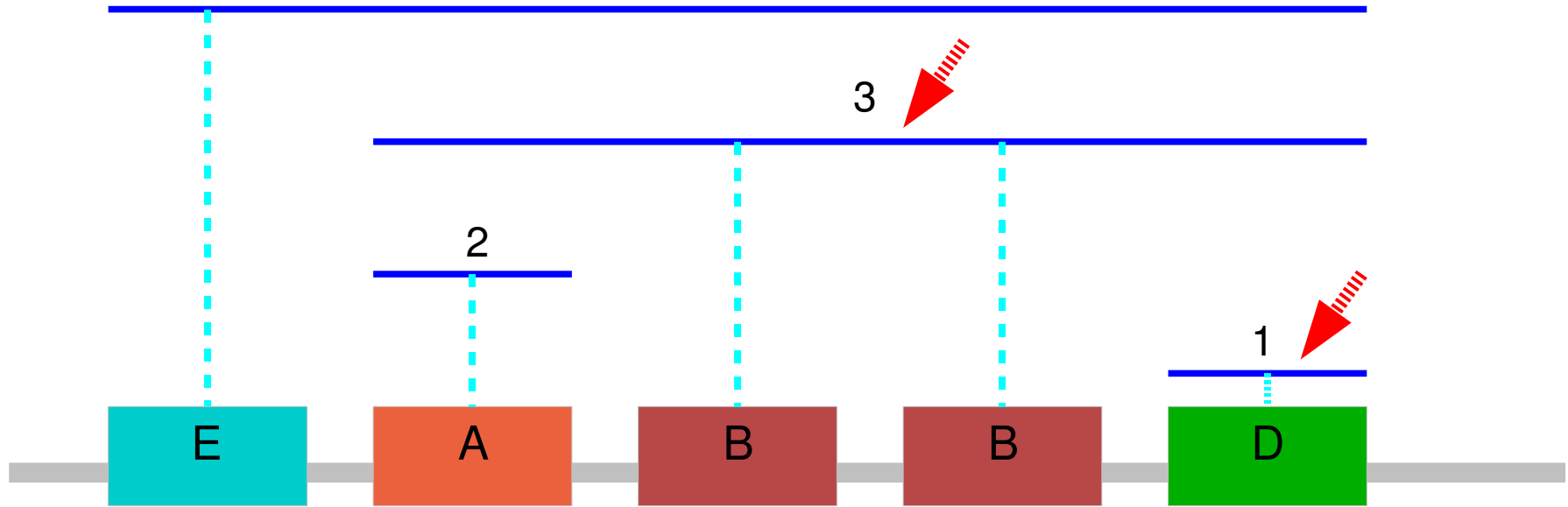
B = 1 C = 2 A = 3

C = 1 A = 2

A = 1

Rank_(Y,1)[B] = 3

Int[k]

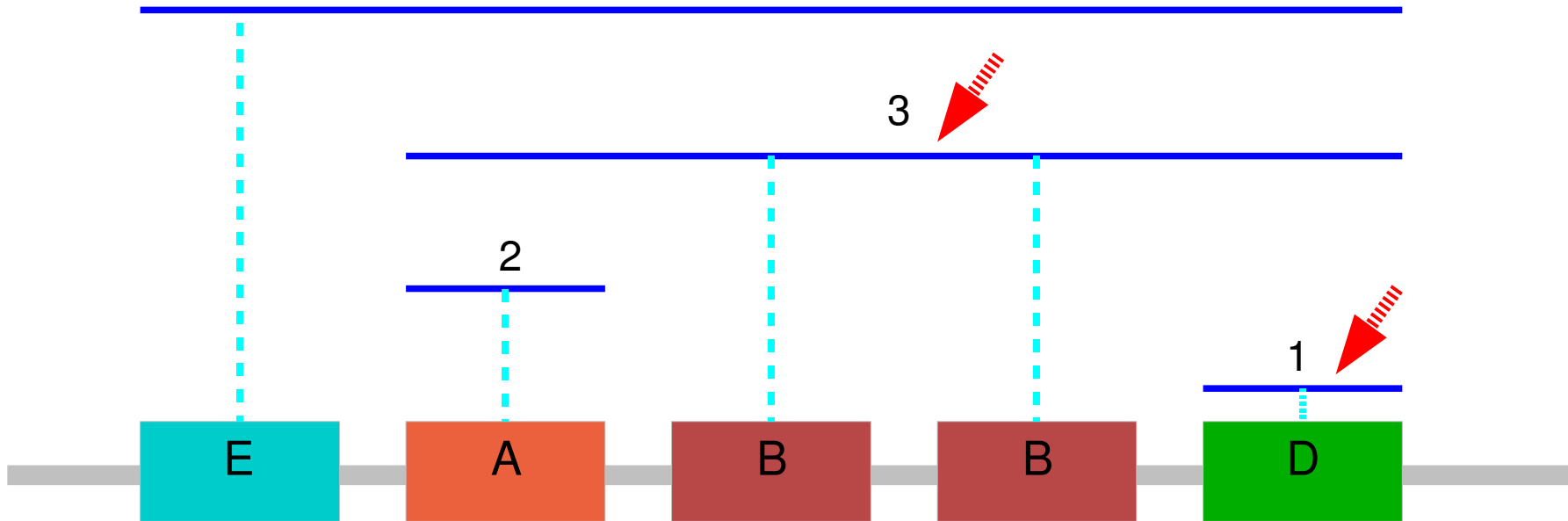


$fo(Y,1) = D A B C$

$B = 1 A = 2 C = 3$

$Rank_{(Y,2)}[A]=2$

Int[k] are nested ! They form a tree. !



$2n$ valid $Int[k]$ at max !

$2nm$ common intervals at maximum

The bound is reached !!

How to identify all them ?

Two approaches

Direct computation (Didier)

$O(nm)$ but

- + Lowest common ancestor (otherwise $O(n m \log n)$)
- + No structure in the output !
- + Complexity does not depend of the input
- + No index

Fingerprint computation on a single string + index+ merge after

- + $O(n+|L_1|\log n + m |L_2| \log m)$ (can be worst than Didier)
- + Structure in the output and possibility of search of fingerprint
- + Complexity does depend of the input
- + Keep the index for further computations

- $S = s_1..s_N$ string of length n
- alphabet Σ of size $|\Sigma|$, not fixed (possibly $O(n)$)

A fingerprint f : set of character(s) of a substring $s_i.. s_j$

General problem:

Compute and represent the set of all fingerprints of S

Examples:

dccbcabbbc

{a} {b} {c} {d} {c,d} {b,c} {a,b} {b,c,d} {a,b,c} {a,b,c,d}

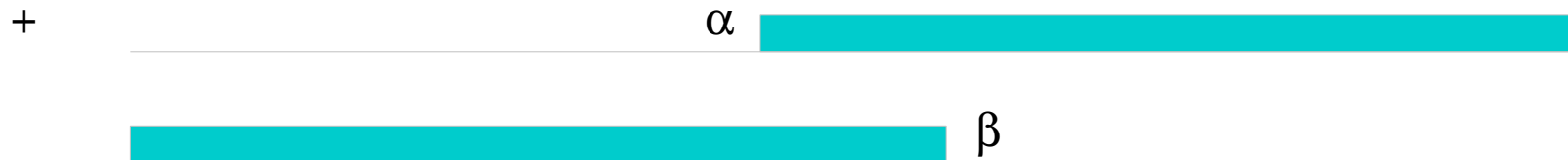
acbdcadad

{a} {b} {c} {d} {a,c} {a,d} {b,c} {b,d} {c,d} {a,b,c} {a,c,d} {b,c,d} {a,b,c,d}

Maximal location $\langle i, j \rangle$ of f



α not in f , β not in f



Number of maximal locations: $L \leq n|\Sigma|$ Complexity of the bound easily reached

But is usually much less

$$\Sigma_k = \{a_1, a_2, \dots, a_k\} \quad w_1 = a_1, w_k = w_{k-1} a_k w_{k-1}$$

$$w_2 = a_1(a_2)a_1, w_3 = (a_1 a_2 a_1) a_3 (a_1 a_2 a_3), \dots$$

$$|w_k| \cdot |L_k| = k \cdot (2^k - 1) \quad |L_k| = 2^{k+1} - (k+2) \quad \longrightarrow \quad |L_k| = o(|w_k| \cdot |L_k|)$$

Naming technique

$\{a,c,e,f\}$ $\Sigma = \{a,b,c,d,e,f,g,h\}$

[7]							
[5]				[6]			
[2]		[2]		[3]		[4]	
[1]	[0]	[1]	[0]	[1]	[1]	[0]	[0]
a	b	c	d	e	f	g	h

$\log |\Sigma| + 1$

[9] ★							
[5]				[8] ★			
[2]		[2]		[3]		[2] ★	
[1]	[0]	[1]	[0]	[1]	[1]	[1]	[0]
						★	

$\{a,c,e,f,g\}$

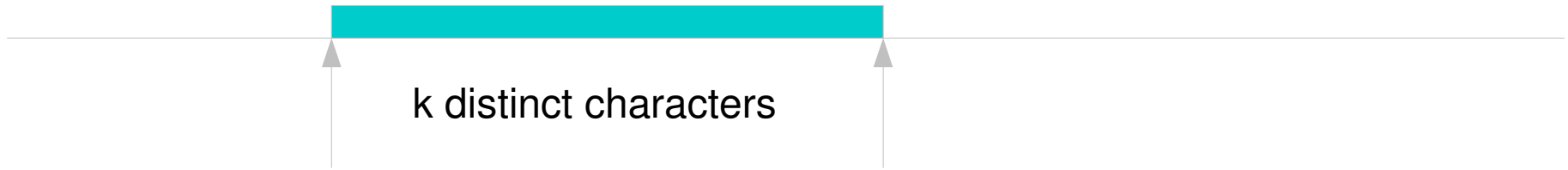
[10] ★							
[5]				[5] ★			
[2]		[2]		[2] ★		[2]	
[1]	[0]	[1]	[0]	[1]	[0]	[1]	[0]
						★	

$\{a,c,e,g\}$

Names = $\{[1],[2],[3],[4],[5],[6],[7],[8],[9],[10]\}$

Fingerprints = $\{[7],[9],[10]\}$

Amir, Apostolico, Landau, Satta 2003



Changing a character: $O(\log |\Sigma| \log n)$ (n new names maximum by level)

One iteration: $n \log |\Sigma| \log n$

Important: different set of names for each iteration

$|\Sigma|$ iterations: $|\Sigma| n \log |\Sigma| \log n$

[4]			
[2]		[3]	
[0]	[0]	[1]	[1]

a b c d

[7]			
[5]		[6]	
[0]	[1]	[1]	[0]

[8]			
[3]		[2]	
[1]	[1]	[0]	[0]

[7]			
[5]		[6]	
[0]	[1]	[1]	[0]

k=2 d c c b c b a b b b c



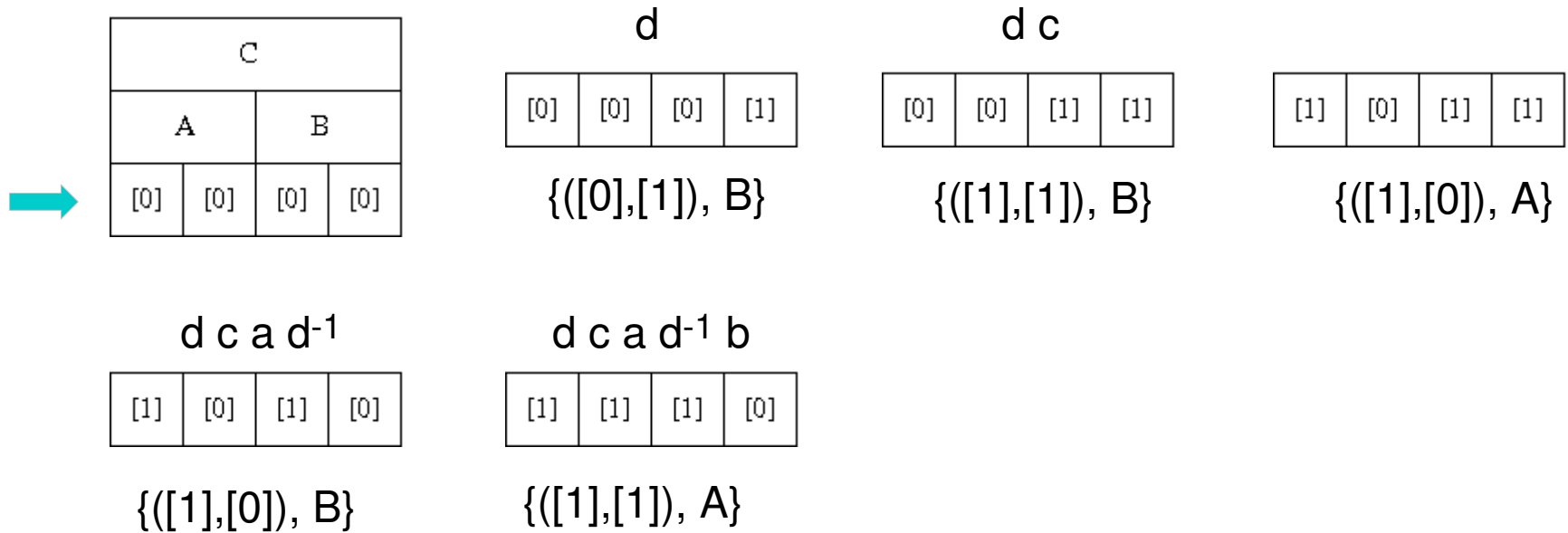
d c c b c b a b b b c



Tsur 2005

List of fingerprints: $d c a d^{-1} b$

$\{d\}, \{c,d\}, \{a,c,d\}, \{a,c\}, \{a,b,c\}$



List of changes:

$\{([0],[0]), A\} \{([0],[0]), B\} \mid \{([0],[1]), B\} \{([1],[1]), B\} \{([1],[0]), A\} \{([1],[0]), B\} \{([1],[1]), A\}$

Radix sort on the pairs + unique -> new names

Tsur 2005

List of changes:

$\{([0],[0]), A\} \{([0],[0]), B\} \mid \{([0],[1]), B\} \{([1],[1]), B\} \{([1],[0]), A\} \{([1],[0]), B\} \{([1],[1]), A\}$

$[2] \rightarrow ([0],[0])$

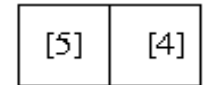
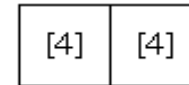
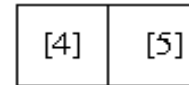
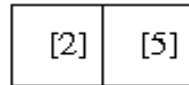
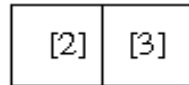
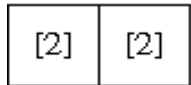
$[3] \rightarrow ([0],[1])$

$[4] \rightarrow ([1],[0])$

$[5] \rightarrow ([1],[1])$

New list:

$\{[2], A\} \{[2], B\} \mid \{[3], B\} \{[5], B\} \{[4], A\} \{[4], B\} \{[5], A\}$



$\{([2],[2]), C\}$

$\{([2],[3]), C\}$

New list: $\{([2],[2]), C\} \mid \{([2],[3]), C\} \{([2],[5]), C\} \{([4],[5]), C\} \{([4],[4]), C\} \{([5],[4]), C\}$

Radix sort, ...

Radix sort: $O(n)$ (bounded integers)

One iteration : $n \log |\Sigma|$ No more name search !



$|\Sigma|$ iterations: $|\Sigma| n \log |\Sigma|$

Problems

- does not depend of L
- distinct names at each iteration

Our approach (2006)

Simple sequence: no repeated character

lfo(i) a b a c e a b a c d

lfo(4)=ceab

a b a c e a b a c d

lfo(2) = bace

Concatenate # to the sequence

Bijection L / proper prefixes of lfo(i)

cea a b a c e a b a c d #

bac a b a c e a b a c d #



Compute all lfo(i) of S#

Our approach (2006)

How to calculate all lfo(i) ?

abcbadca

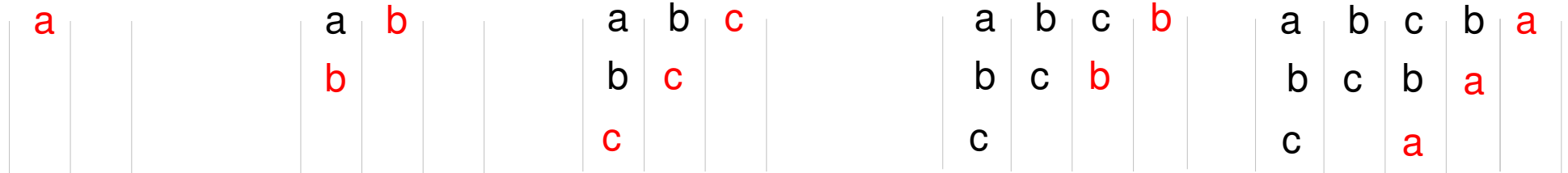
a | bcbadca#

ab | cbadca#

abc | badca#

abcb | adca#

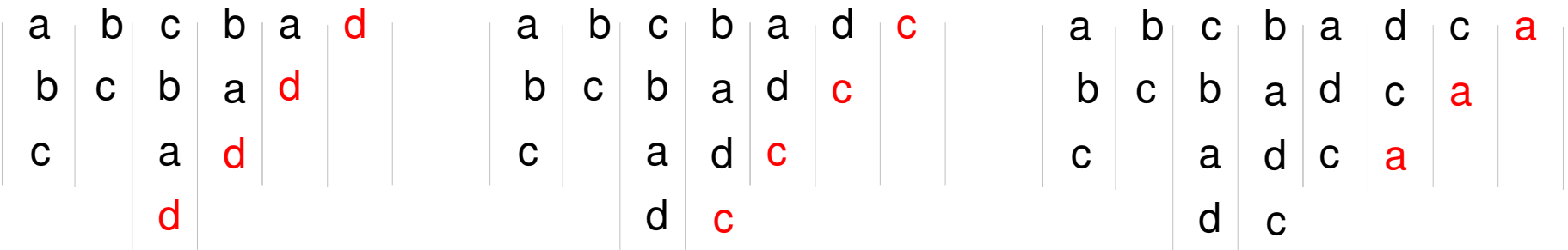
abcba | dca#



abcbad | ca#

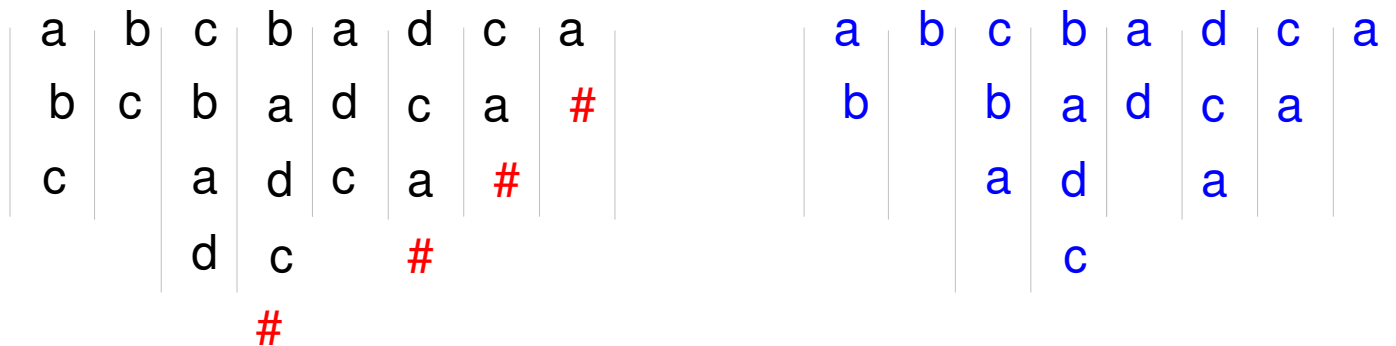
abcbadc | a#

abcbadca | #



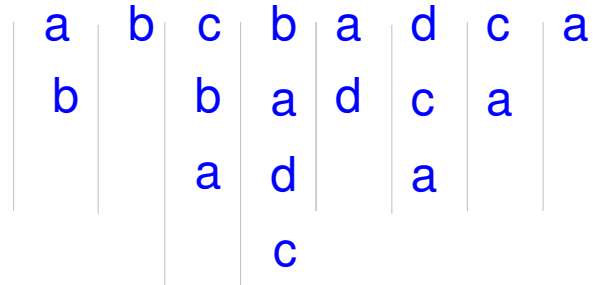
abcbadca#

lfo(i)



Our approach (2006)

Naming all proper prefixes of lfo(i)



n lists:

- Tsur algorithm
- Common names

Simple sequence: $O(|L| \log |\Sigma|)$

General sequence: $O(n + |L| \log |\Sigma|)$

$|L| \leq n |\Sigma|$



Faster or as fast as that of Tsur

- a unique set of names

→ Compute the LCP of two fingerprints in $\log |\Sigma|$

[9] ★							
[5] ★				[8] ★			
[2]		[2]		[3]		[2]	
[1]	[0]	[1]	[0]	[1]	[1]	[1]	[0]

★

[10] ★							
[5] ★				[5] ★			
[2]		[2]		[2] ★		[2]	
[1]	[0]	[1]	[0]	[1]	[0]	[1]	[0]

★

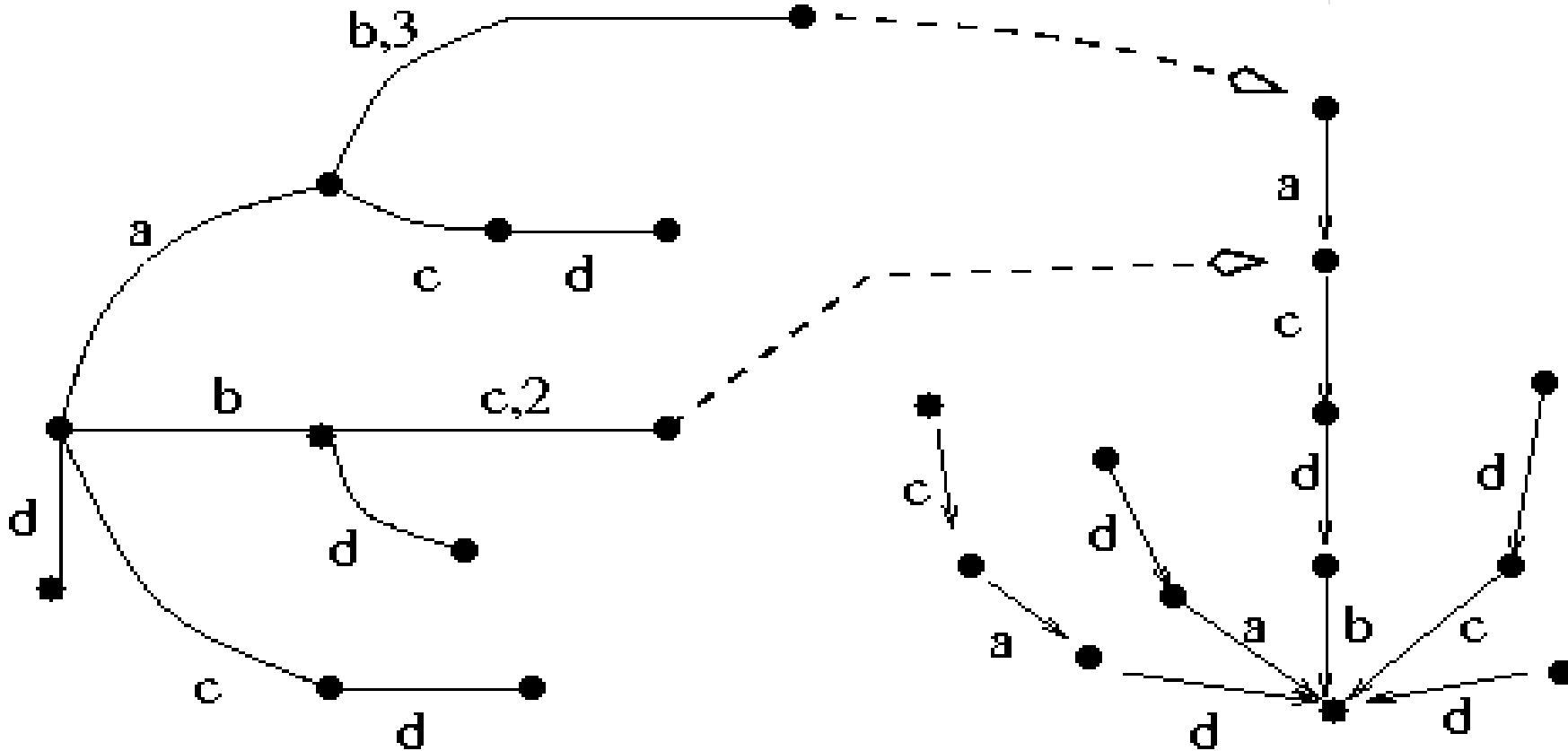
- names sorted by lexicographic order of fingerprints

Fingerprint trie

Chan *et al*, ESA 2007

bdcad

b	d	c	a	d
d	c	a	d	
c		d		
a				



$O(|F|)$ space

$O(|F| \log |\Sigma|)$ time

Search in $O(|f| \log(|f|/|\Sigma|))$

Back to common intervals:


- 1) Build the tree for the first sequence: $O(n+|L_1| \log |\Sigma|)$
- 2) Build the tree for the second sequence: $O(m+|L_2| \log |\Sigma|)$
- 3) Merge the two trees !


Complexity: $O((n+m)+(|L_1|+|L_2|) \log |\Sigma|)$ time.

Open problems

 Memory space reduction

 Order ?

 Approximate fingerprint

 Distance by fingerprints

 2D fingerprints