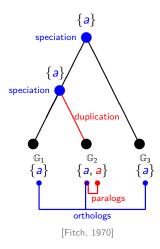
Topics of today:

1. Inferring gene families via family-free DCJ-indel distance

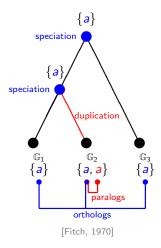
2. Quiz-review

Introduction

gene family: set of genes sharing a common ancestor

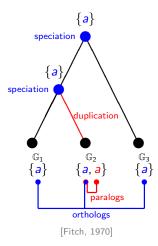


gene family: set of genes sharing a common ancestor



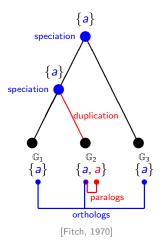
family	\mathbb{G}_1	\mathbb{G}_2	\mathbb{G}_3	
а	1	2	1	ambiguous

gene family: set of genes sharing a common ancestor



# of occurrences											
family \mathbb{G}_1 \mathbb{G}_2 \mathbb{G}_3											
а	1	2	1	ambiguous							
Ь	3	0	1	ambiguous							
с	0	1	1	resolved							
d	1	1	1	resolved							

gene family: set of genes sharing a common ancestor



# of occurrences										
family	\mathbb{G}_1									
а	1	2	1	ambiguous						
b	3	0	1	ambiguous						
с	0	1	1	resolved						
d	1	1	1	resolved						

Identifying the orthologs of two genomes within an ambiguous family is often not so easy

 \downarrow

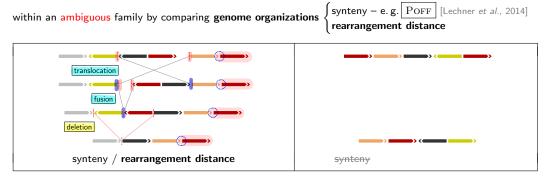
the genome organizations might provide some hints $\ \blacktriangleright$

within an ambiguous family by comparing genome organizations



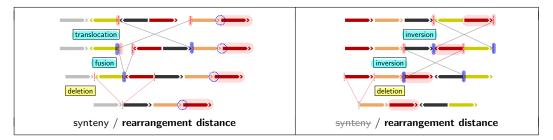
within an ambiguous family by comparing genome organiz	ations { synteny - e.g. Poff [Lechner <i>et al.</i> , 2014]
synteny	

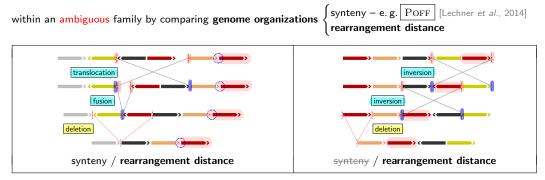
within an ambiguous family by comparing genome organizations synteny – e.g. POFF [Lechner *et al.*, 2014] rearrangement distance



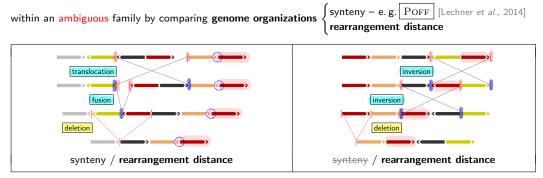
within an ambiguous family by comparing genome organizations \langle

synteny – e. g. POFF [Lechner *et al.*, 2014] rearrangement distance





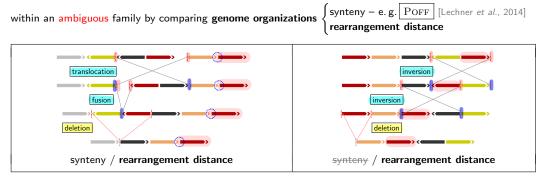
rearrangement distance can identify orthologies that are not "visible" with synteny only [Shi et al., 2010]



rearrangement distance can identify orthologies that are not "visible" with synteny only [Shi et al., 2010]

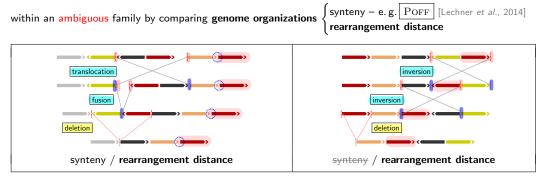
	organizational : change # of chromosomes, positions and orientations of genes
rearrangement types	

types

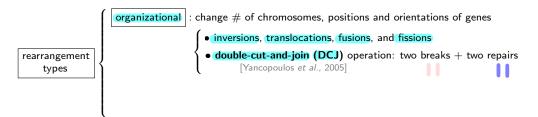


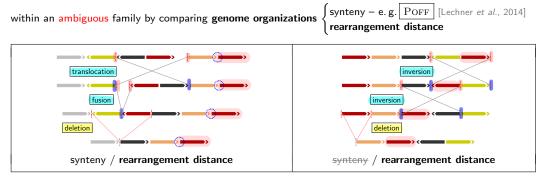
rearrangement distance can identify orthologies that are not "visible" with synteny only [Shi et al., 2010]

organizational : change # of chromosomes, positions and orientations of genes • inversions, translocations, fusions, and fissions rearrangement

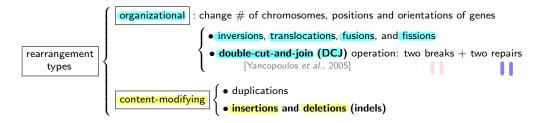


rearrangement distance can identify orthologies that are not "visible" with synteny only [Shi et al., 2010]





rearrangement distance can identify orthologies that are not "visible" with synteny only [Shi et al., 2010]



finitial minimum number of rearrangements transforming one genome into the other

	inversions (+ transl., fus., fis.)	(Hannenhalli and Pevzner, 1995)
only resolved families: polynomial - e.g.	DCJ operations	(Bergeron <i>et al.</i> , 2006)
	DCJ operations and indels	(Braga <i>et al.</i> , 2010)

finitial minimum number for rearrangements transforming one genome into the other

only resolved families: polynomial - e.g. $\begin{cases} inversions (+ transl., fus., fis.) & (Hannenhalli and Pevzner, 1995) \\ DCJ operations & (Bergeron et al., 2006) \\ DCJ operations and indels & (Braga et al., 2010) \end{cases}$

with ambiguous families: NP-hard \rightarrow find a matching that minimizes the distance between the derived genomes with resolved families

finitial minimum number of rearrangements transforming one genome into the other

	inversions (+ transl., fus., fis.)	(Hannenhalli and Pevzner, 1995)
only resolved families: polynomial - e.g.	DCJ operations	(Bergeron et al., 2006)
	DCJ operations and indels	(Braga <i>et al.</i> , 2010)

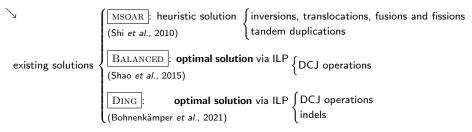
with ambiguous families: NP-hard \rightarrow find a matching that minimizes the distance between the derived genomes with resolved families

\searrow	MSOAR : heuristic solution	\int inversions, translocations, fusions and fissions
	(Shi <i>et al.</i> , 2010)	tandem duplications
existing solutions	{	
	l	

frimum number of rearrangements transforming one genome into the other

	inversions (+ transl., fus., fis.)	(Hannenhalli and Pevzner, 1995)
only resolved families: polynomial - e.g. <	DCJ operations	(Bergeron <i>et al.</i> , 2006)
	DCJ operations and indels	(Braga <i>et al.</i> , 2010)

with ambiguous families: NP-hard \rightarrow find a matching that minimizes the distance between the derived genomes with resolved families



<

<pre>minimum number of rearrangements transforming one genome into the other</pre>	all methods below require the pre-computation of gene families
only resolved families: polynomial - e.	g. dinversions (+ transl., fus., fis.) (Hannenhalli and Pevzner, 1995) DCJ operations (Bergeron <i>et al.</i> , 2006) DCJ operations and indels (Braga <i>et al.</i> , 2010)
with ambiguous families: NP-hard \rightarrow t	find a matching that minimizes the distance between the derived genomes with resolved families
$\sum_{i=1}^{n} \left(\frac{MSOAR}{SOAR} \right) : heu (Shi et al., 2010)$	p) (inversions, translocations, fusions and fissions) (tandem duplications)
existing solutions Shao <i>et al.</i> , 20	optimal solution via ILP {DCJ operations

 $\begin{bmatrix} DING \\ (Bohnenkämper et al., 2021) \end{bmatrix}$ optimal solution via ILP $\begin{cases} DCJ \text{ operations} \\ indels \end{cases}$

Quiz 1

1 Which of the following statements are true?

An ambiguous family occurs more than once in the same genome. A resolved family occurs exactly once in each genome.

2 Two occurrences of the same family in the same genome are called...



3 Genes from a resolved family are called...



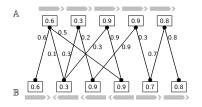
B paralogs

Family-free DCJ-indel distance:

finding orthologs without pre-computed families

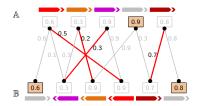
considering { pairwise gene similarities simultaneously genome rearrangements

No family assignment, but pairwise normalized similarities



x=0.1





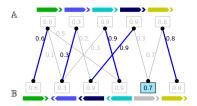
М	$d_{\rm DCJ}^{\rm ID}$	+	M	-	w(M)	+	$w(\widetilde{M})$	=	$wd_{\rm DCJ}^{\rm ID}$
	4	+	4	-	1.7	+	2.3	=	8.6

weighted DCJ-indel distance of M

$$\mathsf{wd}_{ ext{DCJ}}^{ ext{ID}}(\mathbb{A},\mathbb{B},M) = \mathsf{d}_{ ext{DCJ}}^{ ext{ID}}(\mathbb{A},\mathbb{B},M) + |M| - w(M) + w(\widetilde{M})$$

 $M: {\rm matching} \label{eq:matching} {\rm complement} = \widetilde{M}: {\rm set} \mbox{ of } M {\rm -unsaturated} \mbox{ genes}$





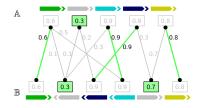
М	$d_{\rm DCJ}^{\rm ID}$	+	M	-	w(M)	+	$w(\widetilde{M})$	=	$wd_{\rm DCJ}^{\rm ID}$
	4	+	4	-	1.7	+	2.3	=	8.6
	3	+	5	_	3.5	+	0.7	=	5.2

weighted DCJ-indel distance of M

$$\mathsf{wd}_{ ext{DCJ}}^{ ext{ID}}(\mathbb{A},\mathbb{B},M) = \mathsf{d}_{ ext{DCJ}}^{ ext{ID}}(\mathbb{A},\mathbb{B},M) + |M| - w(M) + w(\widetilde{M})$$

 $M: {\rm matching} \label{eq:matching} {\rm complement} = \widetilde{M}: {\rm set} \mbox{ of } M {\rm -unsaturated} \mbox{ genes}$

No family assignment, but pairwise normalized similarities



М	$d_{\rm DCJ}^{\rm ID}$	+	M	_	w(M)	+	$w(\widetilde{M})$	=	$wd_{\rm DCJ}^{\rm ID}$
	4	+	4	-	1.7	+	2.3	=	8.6
	3	+	5	_	3.5	+	0.7	=	5.2
	3	+	4	-	3.2	+	1.3	=	5.1

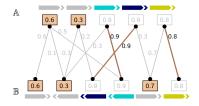
 \square \subset \square

weighted DCJ-indel distance of M

$$\mathsf{wd}_{ ext{DCJ}}^{ ext{ID}}(\mathbb{A},\mathbb{B},M) = \mathsf{d}_{ ext{DCJ}}^{ ext{ID}}(\mathbb{A},\mathbb{B},M) + |M| - w(M) + w(\widetilde{M})$$

$$\label{eq:matching} \begin{split} & M: {\rm matching} \\ {\rm complement} = \widetilde{M}: {\rm set} \mbox{ of } M {\rm -unsaturated genes} \end{split}$$

No family assignment, but pairwise normalized similarities



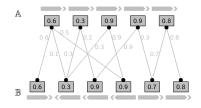
М	$d_{\rm DCJ}^{\rm ID}$	+	M	_	w(M)	+	$w(\widetilde{M})$	=	$wd_{\rm DCJ}^{\rm ID}$
	4	+	4	-	1.7	+	2.3	=	8.6
	3	+	5	-	3.5	+	0.7	=	5.2
	3	+	4	-	3.2	+	1.3	=	5.1
	3	+	3	-	2.6	+	2.5	=	5.9

weighted DCJ-indel distance of M

$$\mathsf{wd}_{\mathrm{DCJ}}^{\mathrm{ID}}(\mathbb{A},\mathbb{B},M) = \mathsf{d}_{\mathrm{DCJ}}^{\mathrm{ID}}(\mathbb{A},\mathbb{B},M) + |M| - w(M) + w(\widetilde{M})$$

 $M: {\rm matching} \label{eq:matching} {\rm complement} = \widetilde{M}: {\rm set} \mbox{ of } M {\rm -unsaturated} \mbox{ genes}$

No family assignment, but pairwise normalized similarities



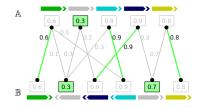
М	$d_{\rm DCJ}^{\rm ID}$	+	M	-	w(M)	+	$w(\widetilde{M})$	=	$wd_{\rm DCJ}^{\rm ID}$
	4	+	4	-	1.7	+	2.3	=	8.6
	3	+	5	-	3.5	+	0.7	=	5.2
	3	+	4	-	3.2	+	1.3	=	5.1
	3	+	3	-	2.6	+	2.5	=	5.9
Ø	2	+	0	-	0	+	7.7	=	9.7

weighted DCJ-indel distance of M

$$\mathsf{wd}_{ ext{DCJ}}^{ ext{ID}}(\mathbb{A},\mathbb{B},M) = \mathsf{d}_{ ext{DCJ}}^{ ext{ID}}(\mathbb{A},\mathbb{B},M) + |M| - w(M) + w(\widetilde{M})$$

$$\label{eq:matching} \begin{split} & M: {\rm matching} \\ {\rm complement} = \widetilde{M}: {\rm set} \mbox{ of } M {\rm -unsaturated genes} \end{split}$$

No family assignment, but pairwise normalized similarities



М	$d_{\rm DCJ}^{\rm ID}$	+	M	_	w(M)	+	$w(\widetilde{M})$	=	$wd_{\rm DCJ}^{\rm ID}$
	4	+	4	-	1.7	+	2.3	=	8.6
	3	+	5	-	3.5	+	0.7	=	5.2
	3	+	4	-	3.2	+	1.3	=	<u>5.1</u>
	3	+	3	-	2.6	+	2.5	=	5.9
Ø	2	+	0	-	0	+	7.7	=	9.7



$$\mathsf{wd}_{\mathrm{DCJ}}^{\mathrm{ID}}(\mathbb{A},\mathbb{B},M) = \mathsf{d}_{\mathrm{DCJ}}^{\mathrm{ID}}(\mathbb{A},\mathbb{B},M) + |M| - w(M) + w(\widetilde{M})$$

$$\label{eq:matching} \begin{split} & M: {\rm matching} \\ {\rm complement} = \widetilde{M}: {\rm set} \mbox{ of } M {\rm -unsaturated genes} \end{split}$$

family-free DCJ-indel distance

$$\mathsf{ffd}^{\mathrm{ID}}_{\mathrm{DCJ}}(\mathbb{A},\mathbb{B}) = \min_{M} \{\mathsf{wd}^{\mathrm{ID}}_{\mathrm{DCJ}}(\mathbb{A},\mathbb{B},M)\}$$

$\mathrm{DIFF}M:$ computing an optimal matching via ILP

Based on the weighted multirelational graph

$$\begin{array}{c|cccc} \textbf{Objective:} & \mbox{Minimize} & \sum_{a \in E_{\xi}} x_a - \sum_{i \leq |V|} \tilde{e}_i + \sum_{k \in K} s_k + \frac{1}{2} \sum_{a \in E} t_a & -\frac{1}{2} \sum_{a \in E_{\xi}} w_a x_a + \sum_{e \in E_{ID}} w_e x_e \\ \hline & & \hline & \hline & & \hline & & \hline & \hline & \hline & & \hline & \hline & & \hline & & \hline & & \hline & & \hline & \hline & \hline & \hline & & \hline & \hline & \hline & \hline & \hline & & \hline & \hline & \hline & \hline & \hline & \hline & & \hline & \hline & \hline & \hline & & \hline & & \hline & \hline & \hline & \hline & & \hline & & \hline & \hline & \hline & \hline & & \hline & \hline & \hline & \hline & & \hline & \hline & \hline & \hline & & \hline & \hline$$

weights appear only in the objective function

Domains:		
(D.01)	$x_a \in \{0, 1\}$	$\forall a \in E$
(D.02)	$0 \ \leq \ell_{i} \leq i$	$\forall \ 1 \leq i \leq V $
(D.03)	$\tilde{e}_{j} \in \{0,1\}$	$\forall \ 1 \leq i \leq V $
(D.04)	$r_V \;\in\; \{0,1\}$	$\forall v \in V$
(D.05)	$t_{a} \in \{0, 1\}$	$\forall a \in E$
(D.06)	$s_k \in \{0,1\}$	$\forall k \in K$

(Rubert et al., 2021)

(Bohnenkämper *et al.*, 2021) (Shao *et al.*, 2015)

Generation of gene families

Generating gene families via genome rearrangements

Inversion of the typical setting gene families \Rightarrow genome rearrangements

Generating gene families via genome rearrangements

Inversion of the typical setting gene families \leftarrow genome rearrangements

Generating gene families via genome rearrangements

Inversion of the typical setting gene families <= genome rearrangements

The optimal matchings given by the ILP | DIFFM | are the basis

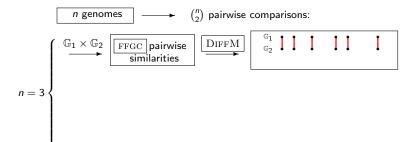
for a generator of gene families of a set of genomes

∜ DIFFMGC

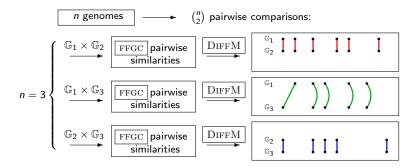
Pipeline of DIFFMGC: FFGC + DIFFM + integration

n genomes $(n \\ 2)$ pairwise comparisons:

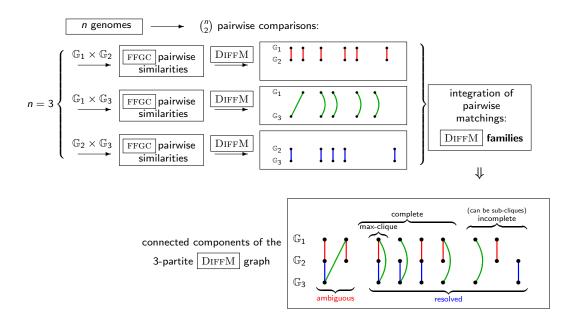
Pipeline of DIFFMGC: FFGC + DIFFM + integration



Pipeline of DIFFMGC: FFGC + DIFFM + integration



Pipeline of DIFFMGC: FFGC + DIFFM + integration



Quiz 2

- 1 The integration of the pairwise matchings given by $\mathrm{D}\mathrm{IFFM}$ is....
 - A a similarity graph
 - B a multipartite graph
 - C a multirelational graph
- 2 Which of the following statements are true?

Any complete resolved family is a max-clique in the DIFFM graph.
 Any max-clique in the DIFFM graph is a complete resolved family.
 Any clique in the DIFFM graph is a resolved family.

3 A max-clique in the DIFFM graph of *n* genomes is composed of...

A
$$n^2$$
 vertices and $\frac{n(n-1)}{2}$ edges

B *n* vertices and
$$\frac{n(n-1)}{2}$$
 edges

C *n* vertices and n^2 edges

mertics

n (m 1) ench edge 2 vortices



PROTEINORTHO : reciprocal best alignment heuristic (Lechner *et al.*, 2011)

PROTEINORTHO : reciprocal best alignment heuristic (Lechner *et al.*, 2011)

POFF : extension of **PROTEINORTHO** incorporating synteny (Lechner *et al.*, 2014)

PROTEINORTHO : reciprocal best alignment heuristic (Lechner *et al.*, 2011)

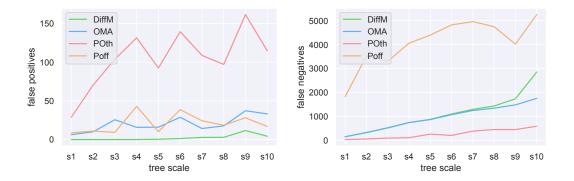
POFF : extension of **PROTEINORTHO** incorporating synteny (Lechner *et al.*, 2014)



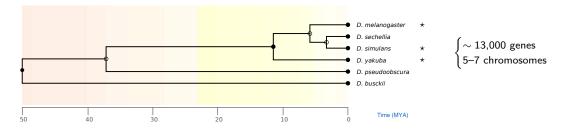
similarities and phylogeny (Dessimoz et al., 2005 & Altenhoff et al., 2019)

Simulated genomes

- ▶ 80 datasets with 10 extant genomes each, with Zombi (Davín et al., 2019)
- High true positive rates for all methods (omitted)

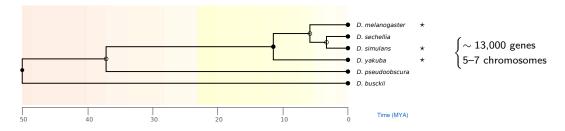


Drosophila genomes



FLYBASE : reference families for *D. melanogaster*, *D. simulans* and *D. yakuba*

Drosophila genomes

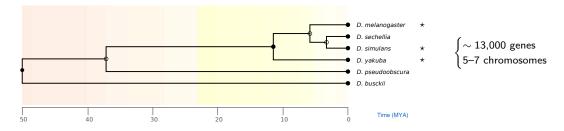


FLYBASE : reference families for *D. melanogaster*, *D. simulans* and *D. yakuba*

		FlyBase families inferred by					
	FlyBase	DIFFMGC	Poff	POrtho	Ома		
resolved	11659	11515	10487	11396	11383		
complete	10809	10713	9745	10594	10594		

all methods found at least 90% of all resolved and complete FLYBASE families

Drosophila genomes



FLYBASE : reference families for *D. melanogaster*, *D. simulans* and *D. yakuba*

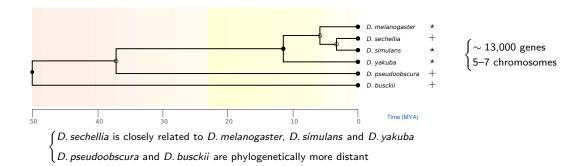
		FlyBase families inferred by					
	FlyBase	DIFFMGC	Poff	POrtho	Oma		
resolved	11659	11515	10487	11396	11383		
complete	10809	10713	9745	10594	10594		

all methods found at least 90% of all resolved and complete FLYBASE families

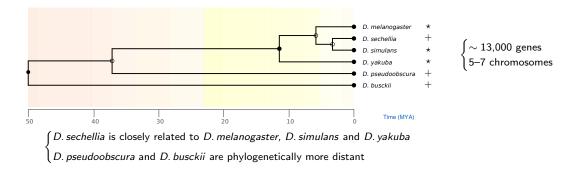
DIFFMGC achieved the highest agreement with FLYBASE

99% of DIFFMGC complete families are max-cliques in the 3-partite DIFFM graph

Results for six Drosophilas



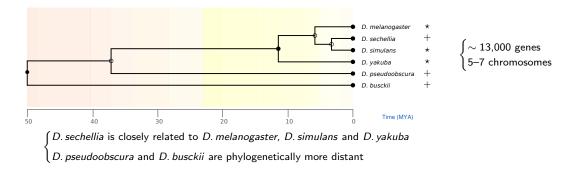
Results for six Drosophilas



Number of families and of resolved and complete families given by the different methods

	DIFFMGC	Poff	POrtho	Oma
total	12885	13282	12746	12660
resolved	12549	13050	11844	11848
complete	8010	8894	8429	8387

Results for six Drosophilas

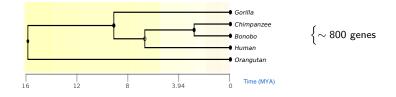


Number of families and of resolved and complete families given by the different methods

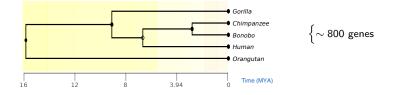
	DIFFMGC	Poff	POrtho	Oma
total	12885	13282	12746	12660
resolved	12549	13050	11844	11848
complete	8010	8894	8429	8387

80% of $\boxed{\mathrm{DIFFMGC}}$ complete families are max-cliques in the 6-partite $\boxed{\mathrm{DIFFM}}$ graph

X chromosome of five primates



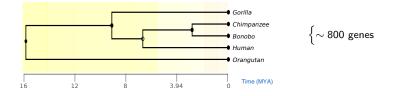
X chromosome of five primates



Number of families and of resolved and complete families given by the different methods:

	DIFFMGC	Poff	POrtho	Oma
total	822	845	782	820
resolved	784	823	715	766
complete	623	609	564	581

X chromosome of five primates



Number of families and of resolved and complete families given by the different methods:

	DIFFMGC	Poff	POrtho	Oma
total	822	845	782	820
resolved	784	823	715	766
complete	623	609	564	581

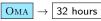
98% of DIFFMGC complete families are max-cliques in the 5-partite DIFFM graph

Running times for six Drosophilas

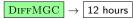
Environment: eight 3GHz cores

pairwise similarities via DIAMOND¹

¹(Buchfink *et al.*, 2015)



pairwise similarities via Smith-Waterman alignments



pairwise similarities via BLAST² (9 hours)

²(Altschul *et al.*, 1990)

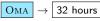
Running times for six Drosophilas

Environment: eight 3GHz cores

 $\fbox{ProteinOrtho} \text{ and } \fbox{Poff} \rightarrow \fbox{5 minutes}$

pairwise similarities via DIAMOND¹

¹(Buchfink *et al.*, 2015)



► DIFFM

pairwise similarities via Smith-Waterman alignments



pairwise similarities via BLAST² (9 hours)

²(Altschul et al., 1990)

	13 took less than 5 minutes				
ILP computations (CPLEX): \sim 3 hours (one took 40 minutes				
	one reached the time limit of 2 hours				
	(with a very small opt. gap of 0.25%)				

References

The potential of family-free rearrangements towards gene orthology inference

(Diego P. Rubert, Daniel Doerr and Marília D. V. Braga)

Journal of Bioinformatics and Computational Biology, vol. 19, No. 06, 2140014 (2021)

Review / Quiz

DCJ-indel model - Path recombinations

With respect to the endpoints:

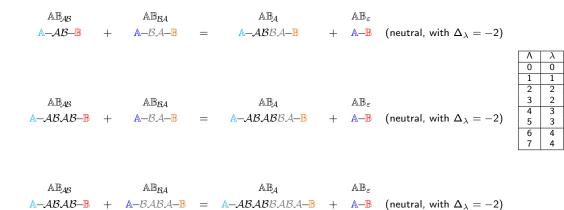
With respect to the runs:

$$A-A + B-B \begin{cases} A-B + A-B & (gaining) \\ A-B + A-B & (gaining) \end{cases} \qquad AB + AB \begin{cases} AA + BB & (\Delta_{\lambda} = -2) \\ ABBA + \varepsilon & (\Delta_{\lambda} = -2) \end{cases}$$
$$A-B + A-B & (gaining) \end{cases} \qquad AB + AB \begin{cases} AA + BB & (\Delta_{\lambda} = -2) \\ ABBA + \varepsilon & (\Delta_{\lambda} = -2) \end{cases}$$
$$A(B) + A \begin{cases} AA + (B) & (\Delta_{\lambda} = -1) \\ AA(B) + \varepsilon & (\Delta_{\lambda} = -1) \end{cases}$$
$$A(B) + A \begin{cases} AA + (B) & (\Delta_{\lambda} = -1) \\ AA(B) + \varepsilon & (\Delta_{\lambda} = -1) \end{cases}$$
$$A(A) + A = \begin{cases} A-A + A-A & (neutral) \\ A-A + A-A & (neutral) \end{cases} \qquad (A)B + B \begin{cases} (A) + BB & (\Delta_{\lambda} = -1) \\ (A)BB + \varepsilon & (\Delta_{\lambda} = -1) \end{cases}$$

 $\label{eq:Deducting path recombinations:} \begin{cases} \mbox{gaining with } \Delta_\lambda = -2 \\ \mbox{gaining with } \Delta_\lambda = -1 \\ \mbox{neutral with } \Delta_\lambda = -2 \end{cases}$

DCJ-indel model - Path recombinations

Putting together (examples):



DCJ-indel model - Path recombinations

Putting together (examples):

$$\begin{array}{rcl} \mathbb{A}\mathbb{B}_{\mathcal{A}\mathcal{B}} & \mathbb{A}\mathbb{B}_{\mathcal{A}\mathcal{B}} & \mathbb{A}\mathbb{B}_{\mathcal{A}\mathcal{B}} & \mathbb{A}\mathbb{B}_{\mathcal{B}} \\ \mathbb{A}-\mathcal{A}\mathcal{B}-\mathbb{B} & + & \mathbb{A}-\mathcal{A}\mathcal{B}-\mathbb{B} & = & \mathbb{A}-\mathcal{A}\mathcal{B}\mathcal{A}\mathcal{B}-\mathbb{B} & + & \mathbb{A}-\mathbb{B} \end{array} (neutral, with $\Delta_{\lambda}=-1$)

$$\begin{array}{rcl} \mathbb{A}\mathbb{B}_{\mathcal{A}\mathcal{B}} & \mathbb{A}\mathbb{B}_{\mathcal{A}\mathcal{B}} & \mathbb{A}\mathbb{B}_{\mathcal{A}\mathcal{B}} & \mathbb{A}\mathbb{B}_{\mathcal{E}} \\ \mathbb{A}-\mathcal{A}\mathcal{B}-\mathbb{B} & + & \mathbb{A}-\mathcal{A}\mathcal{B}-\mathbb{B} & = & \mathbb{A}-\mathcal{A}\mathcal{B}\mathcal{A}\mathcal{B}-\mathbb{B} & + & \mathbb{A}-\mathbb{B} \end{array} (neutral, with $\Delta_{\lambda}=-1$)

$$\begin{array}{rcl} \mathbb{A}\mathbb{B}_{\mathcal{A}\mathcal{B}} & \mathbb{A}\mathbb{B}_{\mathcal{B}\mathcal{A}} & \mathbb{A}\mathbb{B}_{\mathcal{B}} \\ \mathbb{A}-\mathcal{A}\mathcal{B}-\mathbb{B} & + & \mathbb{A}-\mathcal{B}\mathcal{A}-\mathbb{B} & = & \mathbb{A}-\mathcal{A}\mathcal{B}\mathcal{B}\mathcal{A}-\mathbb{B} & + & \mathbb{A}-\mathbb{B} \end{array} (neutral, with $\Delta_{\lambda}=-2$)

$$\begin{array}{rcl} \mathbb{A}\mathbb{B}_{\mathcal{A}\mathcal{B}} & \mathbb{A}\mathbb{B}_{\mathcal{B}\mathcal{A}} & \mathbb{A}\mathbb{B}_{\mathcal{B}} \\ \mathbb{A}-\mathcal{A}\mathcal{B}-\mathbb{B} & + & \mathbb{A}-\mathcal{B}\mathcal{A}-\mathbb{B} & = & \mathbb{A}-\mathcal{A}\mathcal{A}\mathcal{A}-\mathbb{B} \end{array} + & \mathbb{A}-\mathcal{B}\mathcal{B}-\mathbb{B} \end{array} (neutral, with $\Delta_{\lambda}=-2$)

$$\begin{array}{rcl} \mathbb{A}\mathbb{B}_{\mathcal{A}\mathcal{B}} & \mathbb{A}\mathbb{B}_{\mathcal{B}\mathcal{A}} \\ \mathbb{A}-\mathcal{A}\mathcal{B}-\mathbb{B} & + & \mathbb{A}-\mathcal{B}\mathcal{A}-\mathbb{B} \end{array} = & \mathbb{A}-\mathcal{B}\mathcal{A}\mathcal{A}\mathcal{B}-\mathbb{B} + & \mathbb{A}-\mathcal{B} \end{array} (neutral, with $\Delta_{\lambda}=-2$)

$$\begin{array}{rcl} \mathbb{A}\mathbb{B}_{\mathcal{A}\mathcal{B}} & \mathbb{A}\mathbb{B}_{\mathcal{B}\mathcal{A}} \\ \mathbb{A}-\mathcal{A}\mathcal{B}-\mathbb{B} & + & \mathbb{A}-\mathcal{B}\mathcal{A}-\mathbb{B} \end{array} = & \mathbb{A}-\mathcal{B}\mathcal{A}-\mathcal{B}\mathbb{B} \end{array} + & \mathbb{A}-\mathbb{B} \end{array} (neutral, with \Delta_{\lambda}=-2)$$$$$$$$$$$$

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

DCJ and DCJ-indel models - Capping

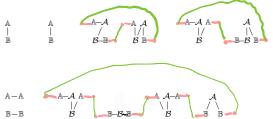
Add caps to close all paths of the graph into cycles, preserving the distance

Canonical capping (no indel edges): maximizes the number of cycles

 $\begin{bmatrix} A \\ B \\ B \end{bmatrix} = \begin{bmatrix} A \\ B \\ B \end{bmatrix} = \begin{bmatrix} A - A \\ B - B \end{bmatrix} = \begin{bmatrix} A - A \\ A \end{bmatrix}$

paths	linking cycle	Δn	Δc	$\Delta(2AB)$	$\Delta_{\rm DCJ}$
AB	(AB)	+0.5	$^{+1}$	-0.5	0
AA + BB	(AA, BB)	$^{+1}$	$^{+1}$	0	0
AA	(AA, Γ _B)	$^{+1}$	$^{+1}$	0	0
BB	$(\mathbb{BB}, \Gamma_{\mathbb{A}})$	$^{+1}$	$^{+1}$	0	0

Singular capping (with indel edges): optimizes the number of cycles and of runs at the same time



paths	linking cycle	Δn	Δc	$\Delta(2\mathbb{AB})$	Δλ	$\Delta_{\mathrm{DCJ}}^{\boldsymbol{\lambda}}$
$AA_{AB} + BB_{AB}$	$(\mathbb{A}\mathbb{A}_{\mathcal{A}\mathcal{B}}, \mathbb{B}\mathbb{B}_{\mathcal{B}\mathcal{A}})$	$^{+1}$	$^{+1}$	0	-2	-2
$2 \times \mathbb{AA}_{AB} + \mathbb{BB}_{A} + \mathbb{BB}_{B}$	$(\mathbb{AA}_{\mathcal{AB}}, \mathbb{BB}_{\mathcal{B}}, \mathbb{AA}_{\mathcal{BA}}, \mathbb{BB}_{\mathcal{A}})$	+2	$^{+1}$	0	-4	-3
$\mathbb{AB}_{\mathcal{AB}} + \mathbb{AB}_{\mathcal{BA}}$	$(\mathbb{AB}_{\mathcal{AB}},\mathbb{AB}_{\mathcal{BA}})$	$^{+1}$	$^{+1}$	-1	-2	$^{-1}$
AB	(AB)	+0.5	$^{+1}$	-0.5	0	0
AA + BB	(AA, BB)	$^{+1}$	$^{+1}$	0	0	0

DCJ-indel model 1

- 1 The indel-potential is defined as...
 - A the number of runs in a component
 - B the smallest number of runs that can be obtained after sorting with internal gaining DCJs
 - C the number of indel-edges in a component
- 2 The indel-potential of a component depends on..
 - \land its number of runs
 - B its number of indel-edges
 - C its length
- 3 The number of runs in a cycle can be...
 - A 0,1,2,4,6,8,...
 - B any non-negative integer
 - C any positive integer

- 4 The number of runs in a path can be...
 - A 0,1,2,4,6,8,...



C any positive integer

DCJ-indel model 2

- 1 A recombination can reduce the overall number of runs by at most...
 - A 1 B 2 C 3
- 2 A recombination can reduce the overall overall indel-potential by at most...

17

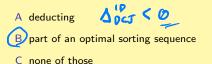
A 1 (B 2) C 3

- 3 A recombination involving a cycle is...
 - A gaining

B neutral

C losing

4 A recombination involving a cycle can be...





DOL5=+2

()

Review-Quiz 1: Inversion model

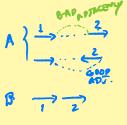
1 In the DCJ model any target adjacency can be reconstructed with an optimal sorting step, but the same is not true for the inversion model because...

🗛 a target adjacency can be bad

- ${\sf B}\,$ a target adjacency can be already present in the genome
- reconstructing a target adjacency can be unsafe
- 2 A cycle is bad when...
 - A it cannot be sorted by inversions
 - B it interleaves another bad cycle Cit contains only bad target adjacencies
- 3 Which data structure helps finding safe inversions?
 - A relational diagram



C component tree



- 4 A bad component can be fixed...
 - A with a neutral inversion
 - B with a split inversion
 - C with a safe inversion

Inversion model 2

- 1 Each leaf of the component tree represents...
 - (A) a bad component
 - B a hurdle
 - C a fortress

- 4 Merging two good (or trivial) components...
 - can merge bad components into a good one
 - B creates a new bad component
 - C is never recommended
- 2 The cost of covering a component tree can be expressed in terms of... gong but bronch
 - A the number of bad nodes
 - B the length of the longest traversal of the tree

the number of leaves cost at he heaves

- 3 Fixing a super hurdle with a neutral inversion
 - A is a good strategy
 - creates a new hurdle
 - C destroys a good component