# CENTER GENOME WITH RESPECT TO THE RANK DISTANCE

Priscila Biller    João P. P. Zanetti    João Meidanis
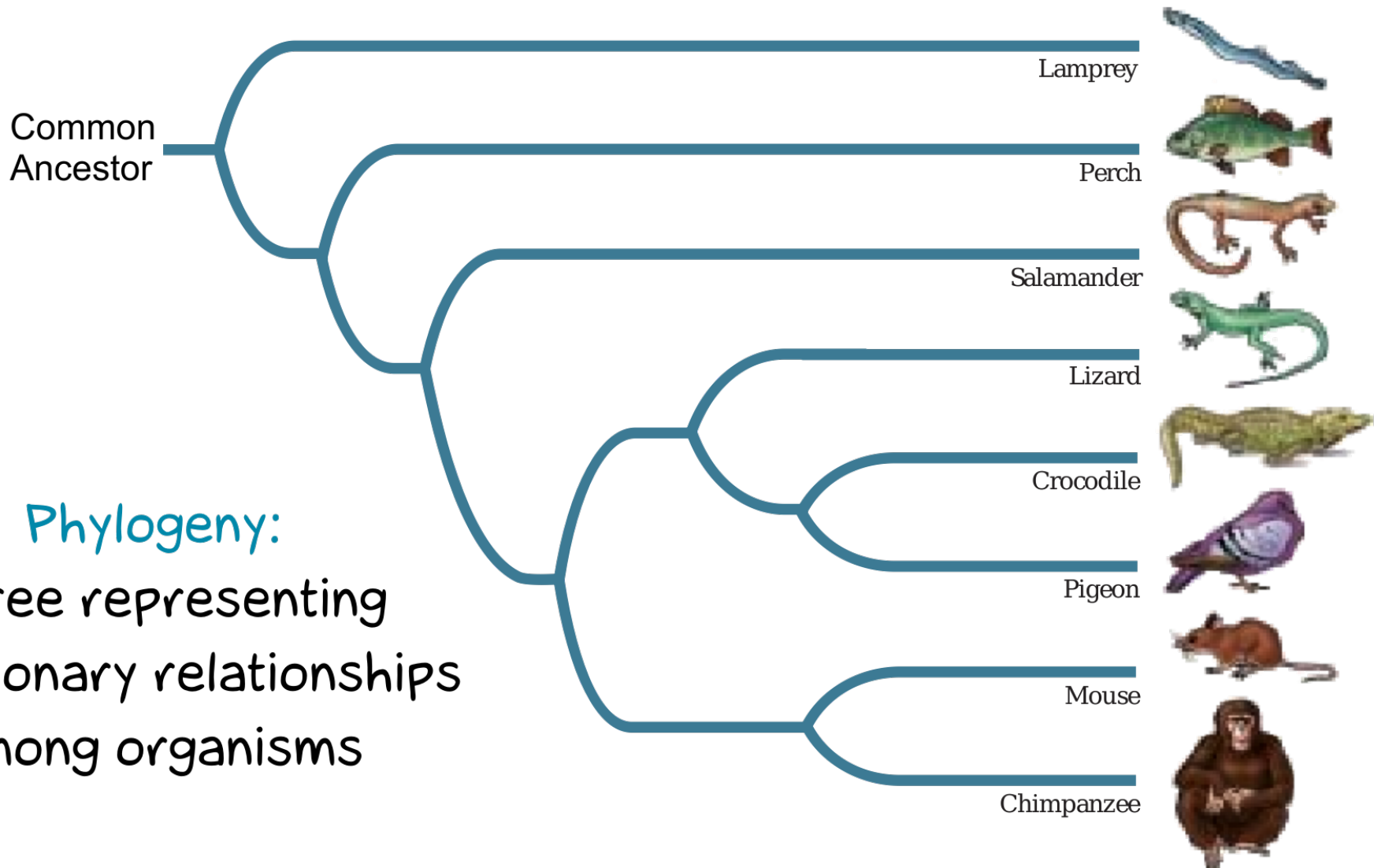
# How to infer ancestral genomes?

Common
Ancestor

Lamprey

Perch

Salamander

Lizard
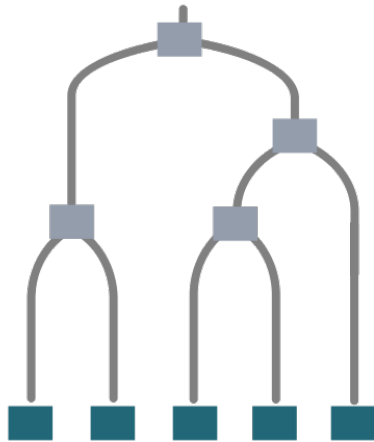
Crocodile

Pigeon

Mouse

Chimpanzee

**Phylogeny:**

a tree representing
evolutionary relationships
among organisms

**Internal nodes:**

ancestors (usually extinct)

**Leaves:**

recent species (known genomes)
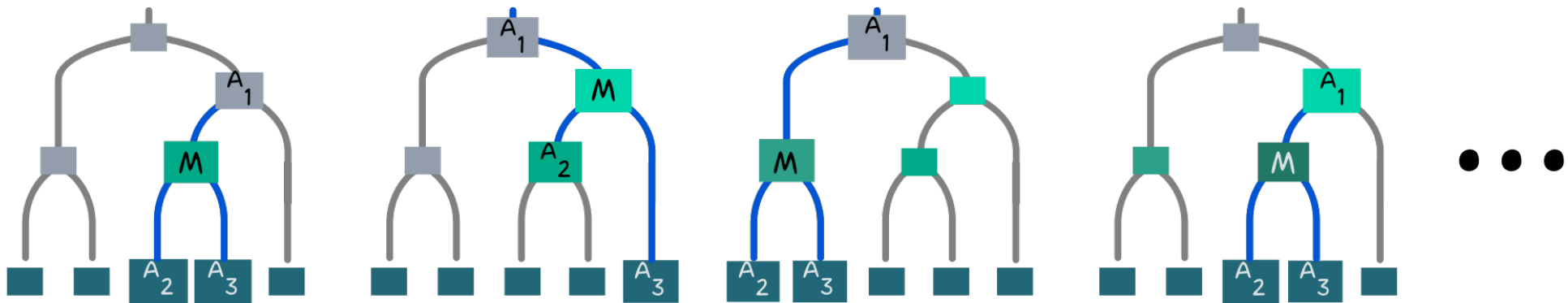
# How to infer ancestral genomes?



■ Known genomes

■ Genomes to be inferred (initially arbitrary)

Usual way to infer ancestors:

Repeatedly compute the median genome M until convergence is reached


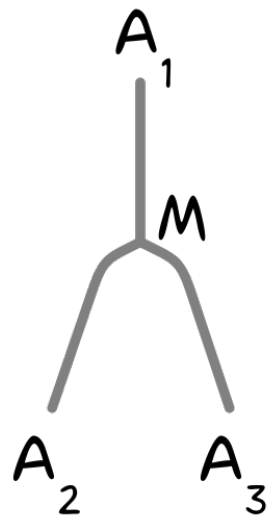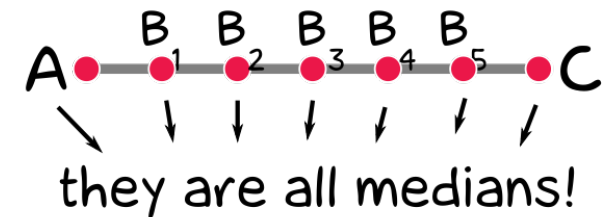
Genome M minimizes the sum of the evolutionary distances:

$$d(A_1, M) + d(A_2, M) + d(A_3, M)$$

# Are median genomes a good way to infer ancestors?

**2 input genomes: solutions are not relevant**

Any genome in an optimal sorting scenario minimizes

$d(A, B) + d(B, C)$, including the input genomes

A •——•$B_1$——•$B_2$——•$B_3$——•$B_4$——•$B_5$——• C

they are all medians!

$A_1$

M

$A_2$ $A_3$

**3 or more input genomes: hard**

NP-hard for most rearrangement

distances (reversal, DCJ, etc.)

# Center genome: an alternative to the median

## Median genome

## Center genome

Input: genomes $A_1$, $A_2$, ...

genomes $A_1$, $A_2$, ...

Goal: find a genome M that minimizes
$d(A_1, M) + d(A_2, M) + d(A_3, M) + ...$

find a genome M that minimizes
$\max(d(A_1, M), d(A_2, M), d(A_3, M), ...)$

2 input genomes:



they are all medians!

❌



more constrained set of solutions

✅

3 or more input genomes:

NP-hard

Open (NP-hard?)

# Ancestral inference: center genomes are an appealing alternative to the median



**Median of 3 genomes** — Hard

Recomputes the same ancestor several times until convergence is achieved

**Center of 2 genomes** — Easy

Direct method

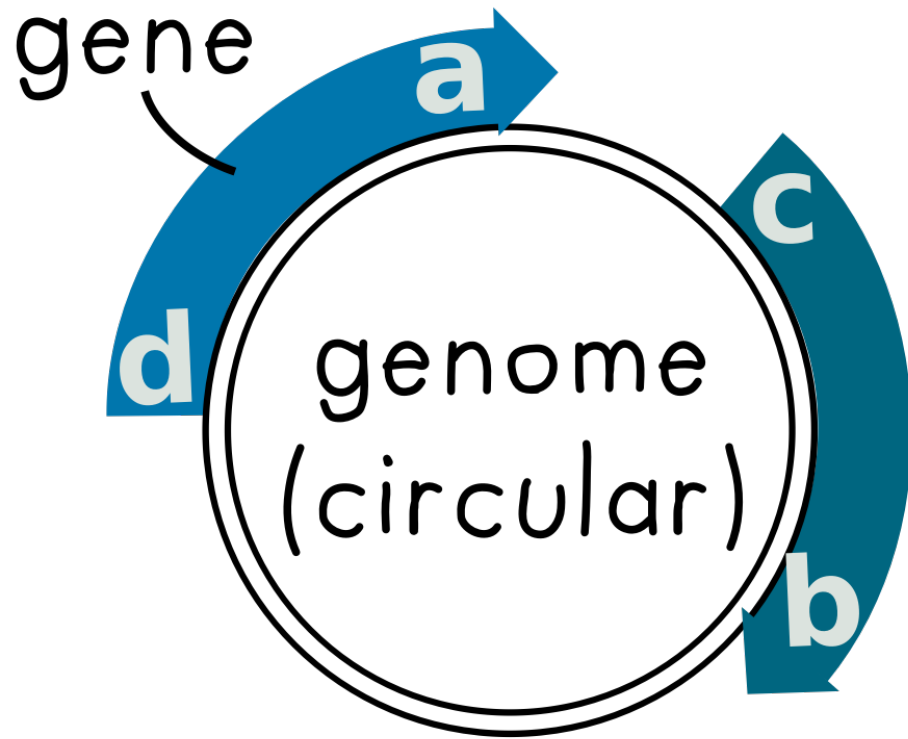# Center genomes with respect to the rank distance

**Rank distance:**

very successfully used in coding

theory since at least 1985

$$d(A, B) = rank(A-B)$$

matrices

Gabidulin Ernst
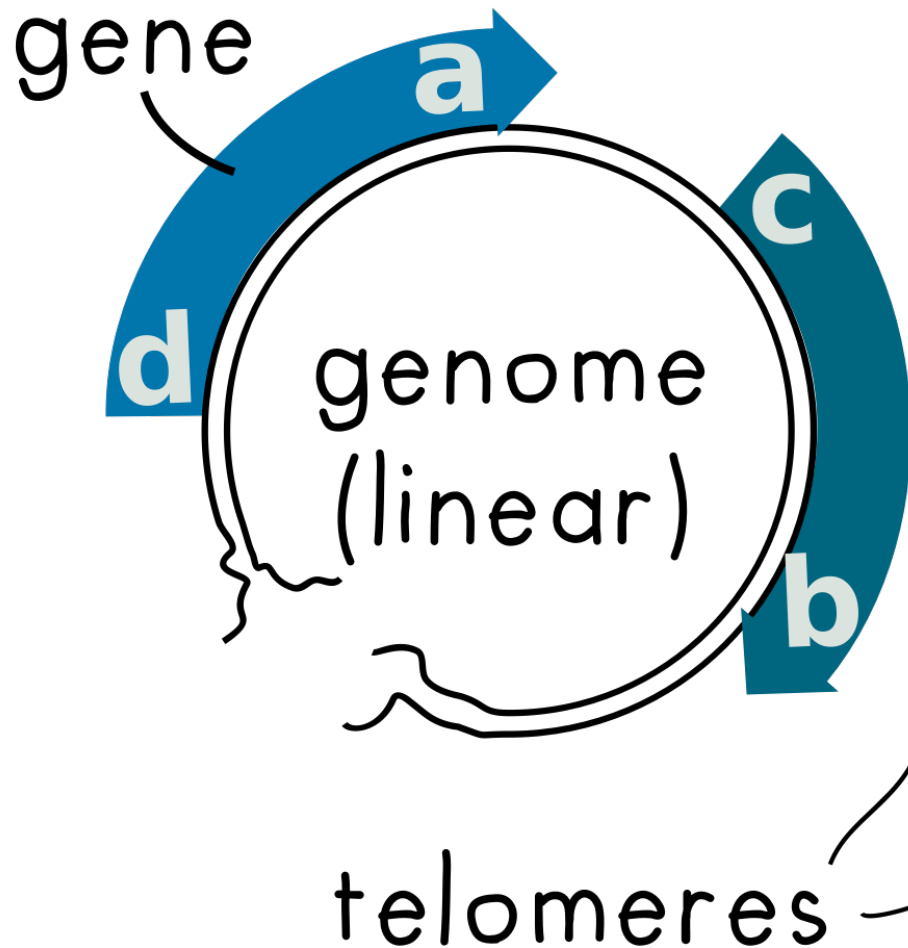
rank(X): dimension of the image

(or column space) of X

# How to represent genomes as matrices?



gene

genome
(circular)

gene extremities

|   | a | b | c | d |
|---|---|---|---|---|
| a | 0 | 0 | 1 | 0 |
| b | 0 | 0 | 0 | 1 |
| c | 1 | 0 | 0 | 0 |
| d | 0 | 1 | 0 | 0 |

# How to represent genomes as matrices?



|   | a | b | c | d |
|---|---|---|---|---|
| a | 0 | 0 | 1 | 0 |
| b | 0 | 1 | 0 | 0 |
| c | 1 | 0 | 0 | 0 |
| d | 0 | 0 | 0 | 1 |

# How to compute the rank distance of two genomes?

$M_1$: 

$M_2$: 

$$d(M_1, M_2) = \mathrm{rank}\left(\begin{array}{cccc} & a & b & c & d \\ a & 0 & 0 & 1 & 0 \\ b & 0 & 0 & 0 & 1 \\ c & 1 & 0 & 0 & 0 \\ d & 0 & 1 & 0 & 0 \end{array} - \begin{array}{cccc} & a & b & c & d \\ a & 0 & 0 & 1 & 0 \\ b & 0 & 1 & 0 & 0 \\ c & 1 & 0 & 0 & 0 \\ d & 0 & 0 & 0 & 1 \end{array}\right)$$

$$= \mathrm{rank}\left(\begin{array}{cccc} & a & b & c & d \\ a & 0 & 0 & 0 & 0 \\ b & 0 & -1 & 0 & 1 \\ c & 0 & 0 & 0 & 0 \\ d & 0 & 1 & 0 & -1 \end{array}\right) = 1$$
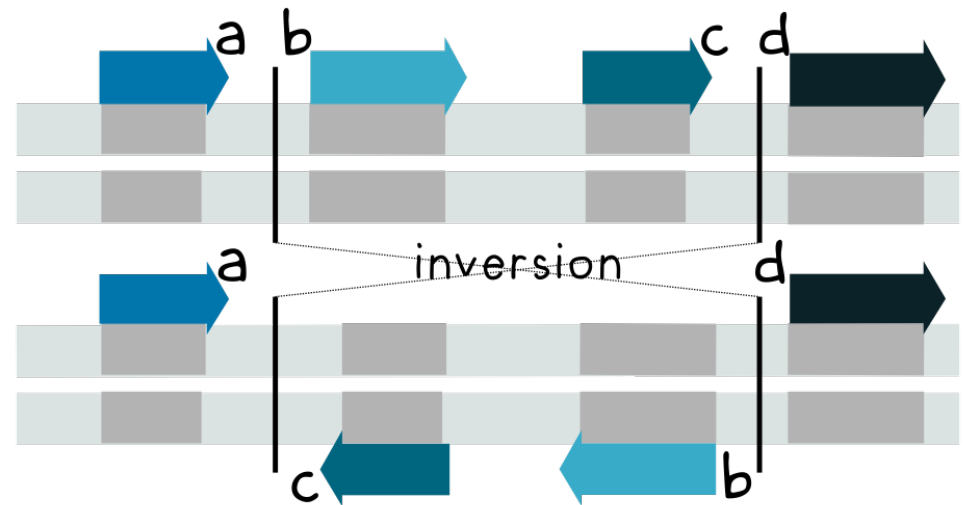
# Applications of the rank distance to genome evolution

Rank distance can also be defined as the minimum number of cuts, joins, and double swaps, with weights 1, 1, and 2
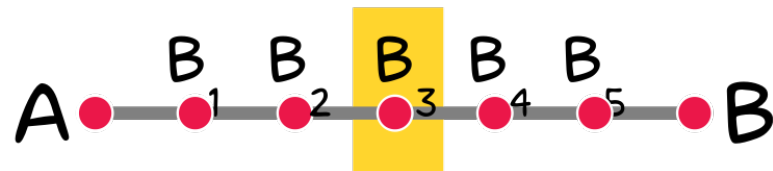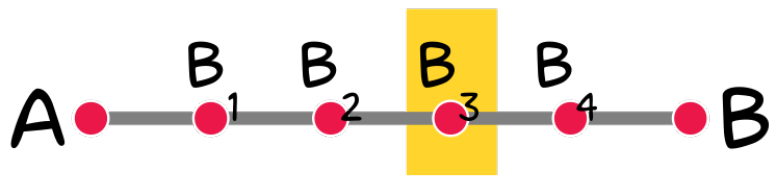


Double swap

$(a, b) \longleftrightarrow (a)$ and $(b)$

$(a, b)$ and $(c, d)$ → $(a, c)$ and $(b, d)$

$(a, b)$ and $(c, d)$ → $(a, d)$ and $(c, b)$
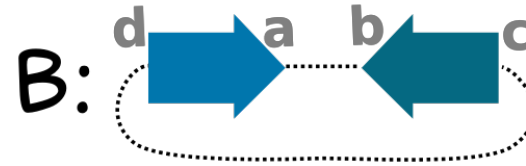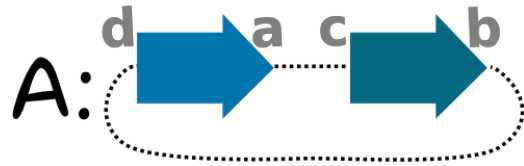
# Center conjecture:

Is there always a genome exactly in the middle?



$$\max(d(A,C), d(B,C)) = \left\lceil \frac{d(A, B)}{2} \right\rceil$$

# Center conjecture (counterexample):

Is there always a genome exactly in the middle?

A:

| | a | b | c | d |
|---|---|---|---|---|
| a | 0 | 0 | 1 | 0 |
| b | 0 | 0 | 0 | 1 |
| c | 1 | 0 | 0 | 0 |
| d | 0 | 1 | 0 | 0 |

B:

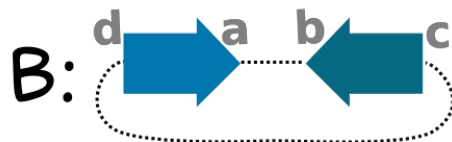| | a | b | c | d |
|---|---|---|---|---|
| a | 0 | 1 | 0 | 0 |
| b | 1 | 0 | 0 | 0 |
| c | 0 | 0 | 0 | 1 |
| d | 0 | 0 | 1 | 0 |

$$d(A, B) = \text{rank} \begin{pmatrix} & a & b & c & d \\ a & 0 & -1 & 1 & 0 \\ b & -1 & 0 & 0 & 1 \\ c & 1 & 0 & 0 & -1 \\ d & 0 & 1 & -1 & 0 \end{pmatrix} = 2$$

# Center conjecture (counterexample):
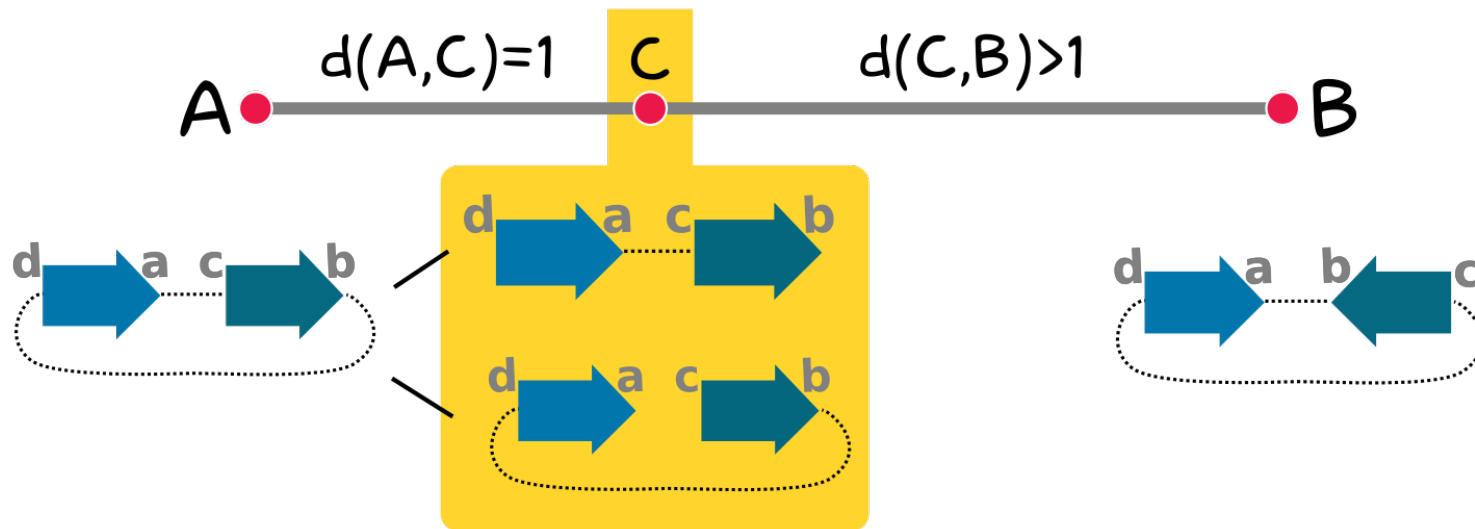
Is there always a genome exactly in the middle?

A:

$$d(A, B) = rank(A-B) = 2$$

B:

There is no genome with distance 1 from both A and B:

$d(A,C)=1$   C   $d(C,B)>1$

A•————————————•————————————•B

In this example:   ~~max(d(A,C), d(B,C)) = $\lceil \dfrac{d(A, B)}{2} \rceil$~~

*Our goal:*

Determine in which cases the center genome

is exactly in the middle and, when not in the middle,

how far it will be from it.



$$\max(d(A,C),\, d(B,C)) = \left\lceil \frac{d(A,\, B)}{2} \right\rceil$$

# Our goal: Where is the center genome?

Given a pair of genomes, there are two cases to consider:

Co-tailed genomes

Not co-tailed genomes

exactly the same telomeres

different telomeres

# Our goal: Where is the center genome?

Co-tailed case (same telomeres)

**Only** double swaps occur in optimal sorting scenarios of co-tailed genomes!

Leonid Chindelevitch et al. (2018)

Double swap (weight = 2)

(a, b) and (c, d)

(a, c) and (b, d)

(a, d) and (c, b)

# Our goal: Where is the center genome?

Co-tailed case (same telomeres)

As a double swap has weight 2, the distance is even!

d(A,B)                    Center genome (C)



0          A = B          LB

2          A C —2— C B          LB + 1

4          A —2— C —2— B          LB

6          A —2— C —2— C —2— B          LB + 1

8          A —2— —2— C —2— —2— B          LB

Lower bound (LB):

$$\max(d(A,C), d(B,C)) = \left\lceil \frac{d(A, B)}{2} \right\rceil$$

# Our goal: Where is the center genome?

## Co-tailed case (same telomeres)
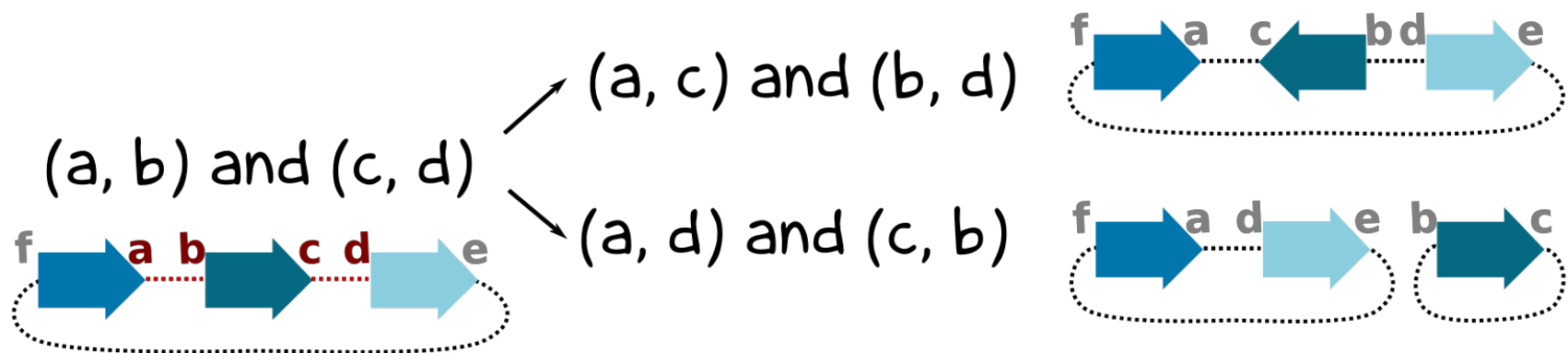
d(A,B) is a multiple of 4 : center genome reaches LB!

Otherwise: LB + 1

✔

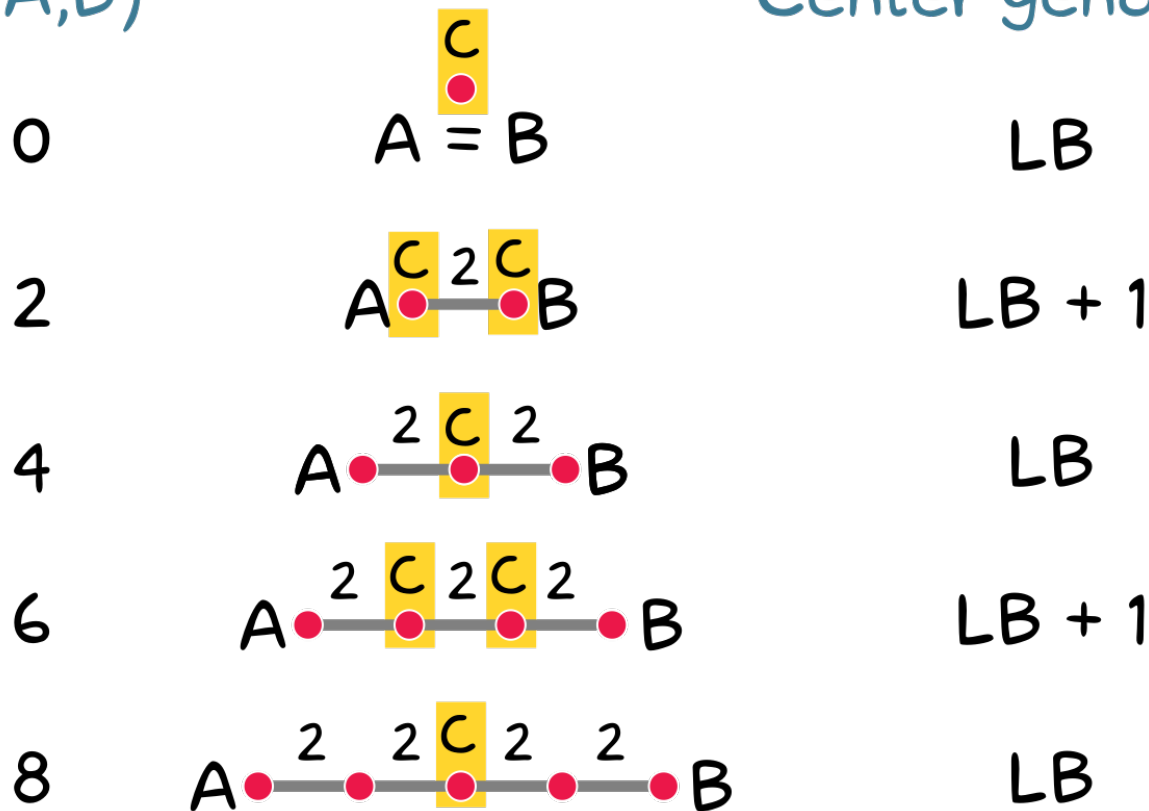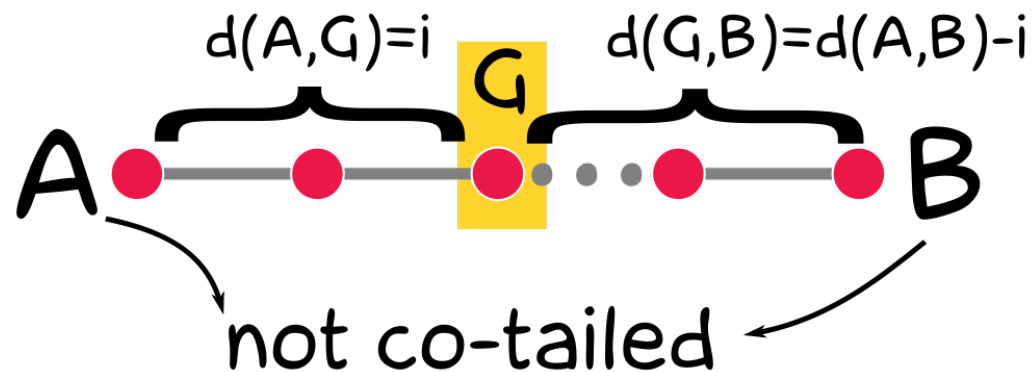| d(A,B) | | Center genome (C) |
|--------|--------|-------------------|
| 0 | A = B | LB |
| 2 | A •—C 2 C—• B | LB + 1 |
| 4 | A •—2 C 2—• B | LB |
| 6 | A •—2 C 2 C 2—• B | LB + 1 |
| 8 | A •—2 2 C 2 2—• B | LB |

# Our goal: Where is the center genome?

## Not co-tailed case (different telomeres)

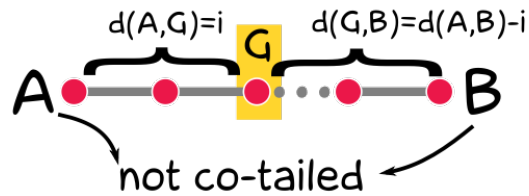At every step $0 \leqslant i \leqslant d(A,B)$ there is a genome G such that

$d(A,G)=i$ and $d(G,B)=d(A,B)-i$...

$d(A,G)=i$    G    $d(G,B)=d(A,B)-i$

A       B

not co-tailed

...so there is a center genome that reaches the LB!

# Our goal: Where is the center genome?

## Not co-tailed case (different telomeres)



$d(A,G)=i$  $d(G,B)=d(A,B)-i$

not co-tailed

If A and B are not co-tailed, then there is a genome at every step between 0 and d(A,B)

**Proof by induction (idea):**  Not co-tailed ➤ It needs a cut!

⋯ Cut an adjacency of A ⋯

also not co-tailed!

A+cut

induction OK!

not co-tailed

✔

co-tailed!

A+cut

induction NOT OK!

not co-tailed

also not co-tailed!

A+double swap

induction OK!

not co-tailed

✔

**Double swap and cut: any order is fine**

cut (e,f)    f  a b  c d  e

ds (a,b),(c,d)    f  a c  b d  e

f  a b  c d  e

ds (a,b),(c,d)    f  a c  b d  e

cut(e,f)

# Summary: Where is the center genome?

Center genome (C)

(1) Co-tailed case (same telomeres)

$$A \overset{2\ \boxed{C}\ 2}{\bullet\!-\!\bullet\!-\!\bullet} B \qquad d(A,B) \text{ is a multiple of 4:} \qquad LB$$

$$A \overset{2\ \boxed{C}\ 2\ \boxed{C}\ 2}{\bullet\!-\!\bullet\!-\!\bullet\!-\!\bullet} B \qquad d(A,B) \text{ is not a multiple of 4:} \qquad LB+1$$
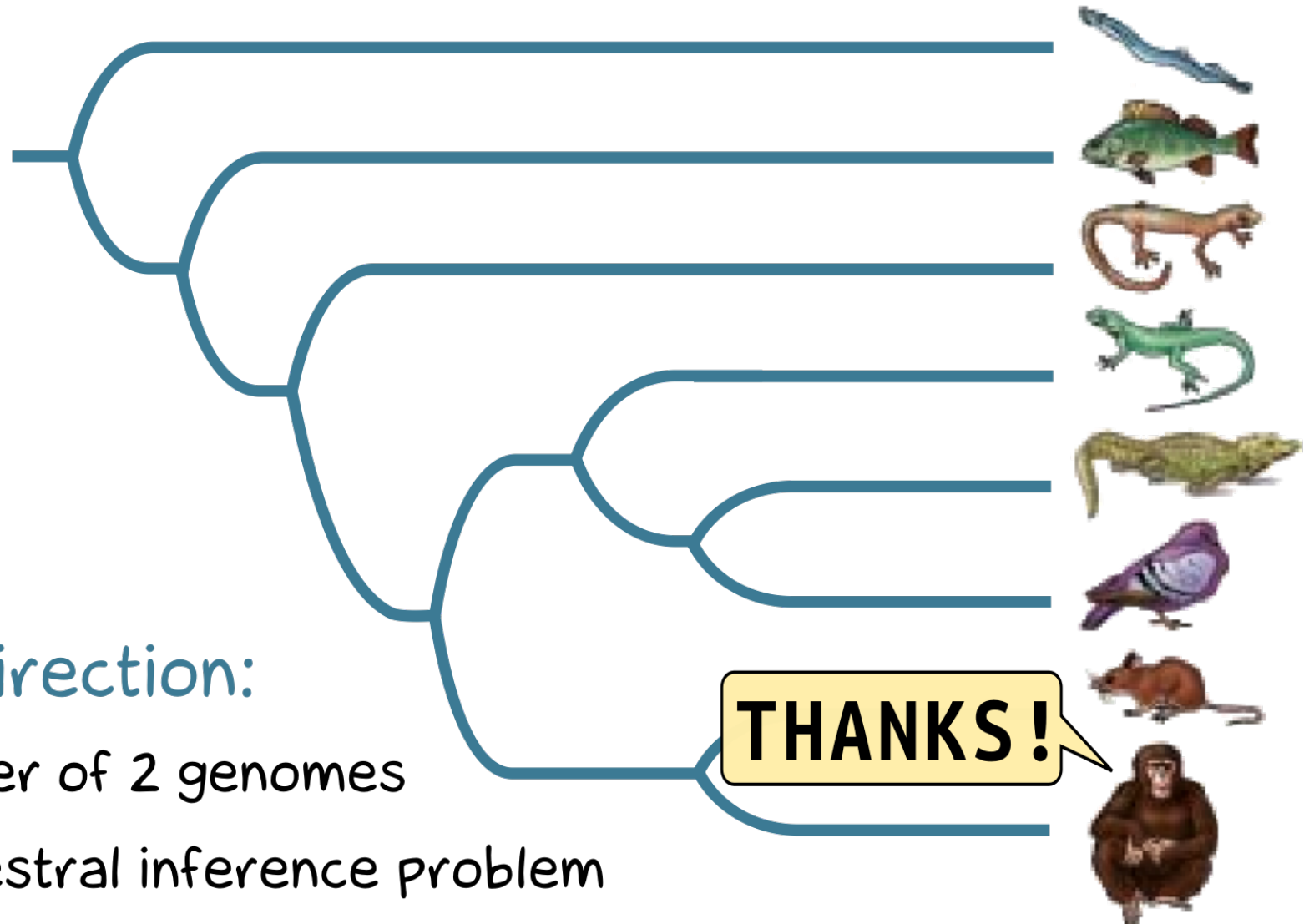
(2) Not co-tailed case (different telomeres)    LB

Lower bound (LB):
$$\max(d(A,C),\ d(B,C)) = \left\lceil \frac{d(A,\ B)}{2} \right\rceil$$

# Center genome with respect to the rank distance



**Future direction:**

Apply center of 2 genomes

to the ancestral inference problem

**THANKS!**