

# Bridging Disparate Views on the DCJ-indel Model

## WABI 2023

Leonard Bohnenkämper

Bielefeld University

September 4, 2023

# The Genomic Distance Problem

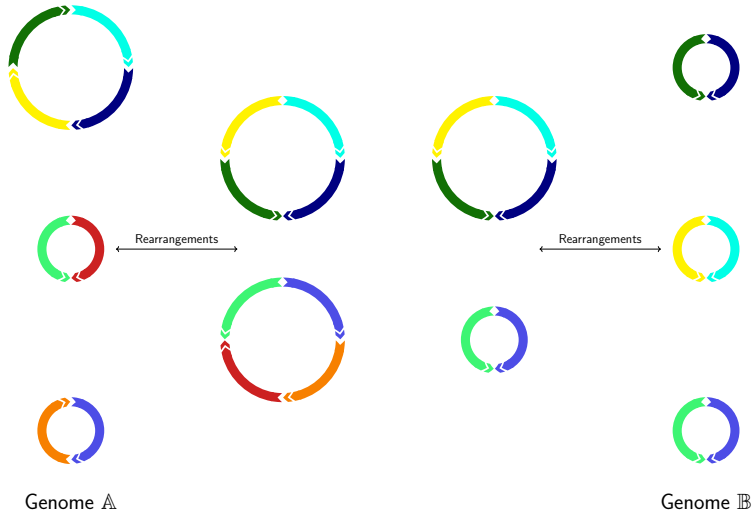


Genome  $\mathbb{A}$

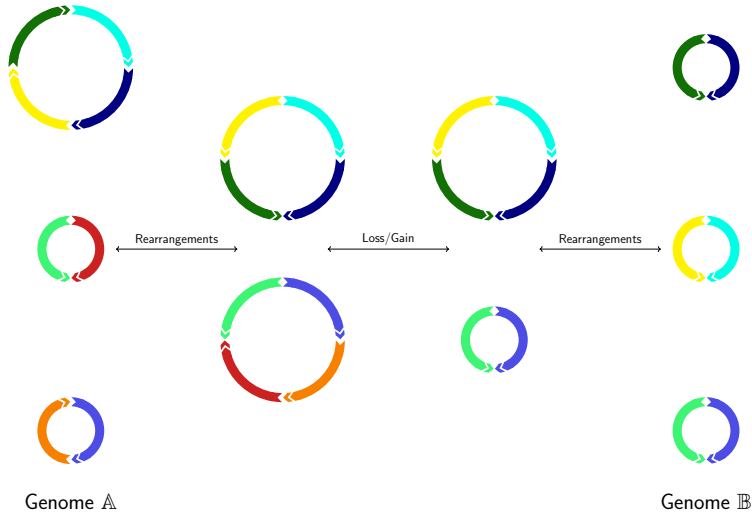


Genome  $\mathbb{B}$

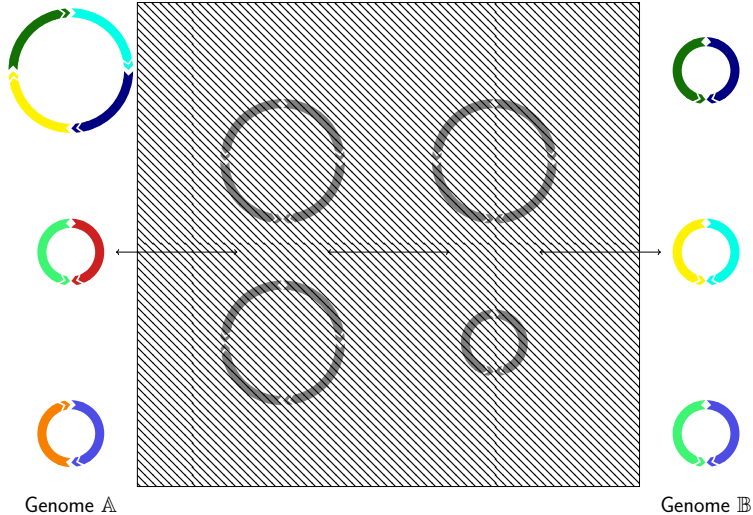
# The Genomic Distance Problem



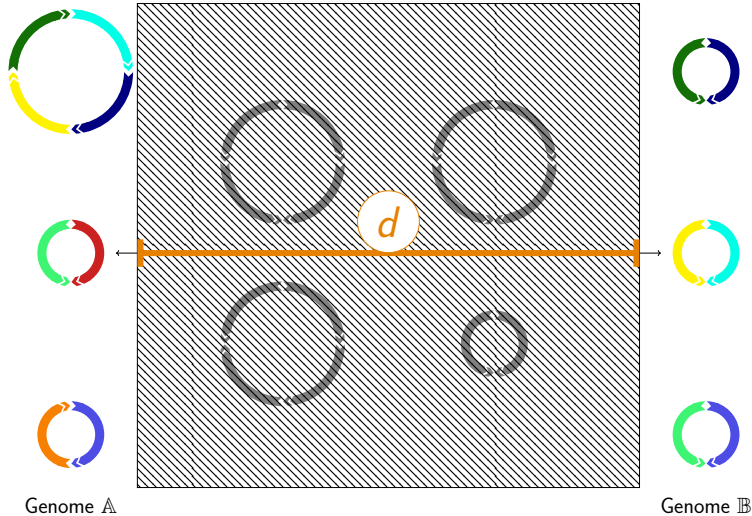
# The Genomic Distance Problem



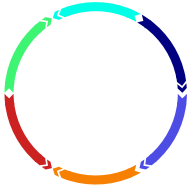
# The Genomic Distance Problem



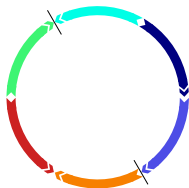
# The Genomic Distance Problem



# Operations: Double-Cut-and-Join (DCJ)

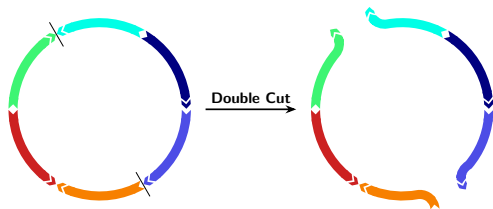


# Operations: Double-Cut-and-Join (DCJ)

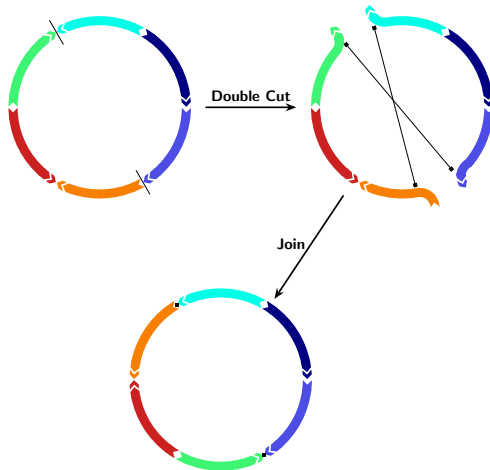




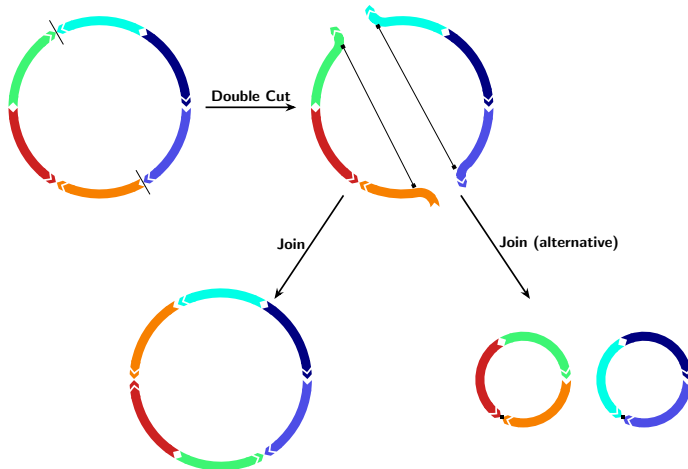
# Operations: Double-Cut-and-Join (DCJ)



# Operations: Double-Cut-and-Join (DCJ)

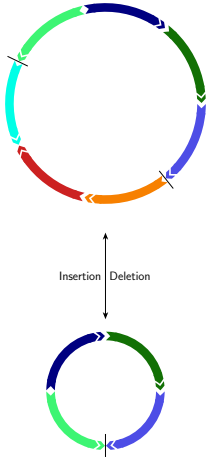


# Operations: Double-Cut-and-Join (DCJ)

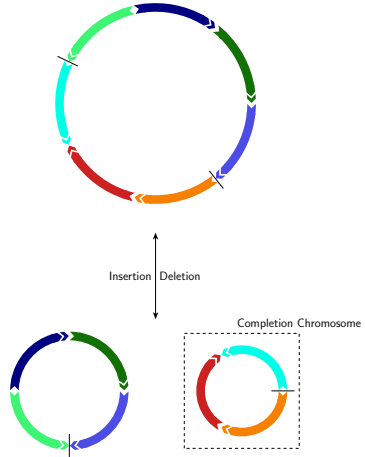


# Operations: Indels

Explicitly (BWS, 2011)



Via Completion (Compeau, 2013)



# Problems with the BWS- and Compeau-Conceptualizations

- ▶ Unwieldy formulas:

# Problems with the BWS- and Compeau-Conceptualizations

- ▶ Unwieldy formulas:

- ▶ BWS:

$$d_{DCJ}^{id}(A, B) = d_{DCJ}(A, B) + \sum_C \lambda(C) \\ - 2P - 3Q - 2T - S - 2M - N$$

# Problems with the BWS- and Compeau-Conceptualizations

- Unwieldy formulas:
- BWS:

$$d_{DCJ}^{id}(A, B) = d_{DCJ}(A, B) + \sum_C \lambda(C) - 2P - 3Q - 2T - S - 2M - N$$

	Sources	Reusable Resul- tants	Unsat. after safe internal ops.					Unsat. after all safe ops. in group					No partner due to
			$A \leq a$	$B \leq b$	$a \leq b$	$A \geq b$	$B \geq a$	$A \leq a$	$B \leq b$	$a \leq b$	$A \geq b$	$B \geq a$	
<i>P</i>													
1	$AA_{AB} + BB_{AB}$		1	1	2	1	1	0	0	0	0	0	
<i>Q</i>													
1	$2AA_{AB} + BB_A + BB_B$		2	2	2	2	2	0	0	0	0	0	
2	$2BB_{AB} + AA_A + AA_B$		2	2	2	2	2	0	0	0	0	0	
<i>T</i>													
1	$AA_{AB} + BB_A + AB_{AB}$		2	2	2	1	1	0	0	0	0	0	
2	$AA_{AB} + BB_B + AB_{AB}$		1	1	2	2	2	0	0	0	0	0	
3	$BB_{AB} + AA_A + AB_{BA}$		2	2	2	1	1	0	0	0	0	0	
4	$BB_{AB} + AA_B + AB_{BA}$		1	1	2	2	2	0	0	0	0	0	
5	$2AA_{AB} + BB_A$	$AA_B$	2	2	2	2	0	0	0	0	2	0	$P_{AB} \rightarrow T1, Q1, P1;$
6	$2AA_{AB} + BB_B$	$AA_A$	2	0	2	2	2	2	0	0	0	0	$P_{AB} \rightarrow T2, Q1, P1;$
7	$2BB_{AB} + AA_A$	$BB_B$	2	2	2	0	2	0	0	0	0	2	$P_{BA} \rightarrow Q2, P1, T3;$
8	$2BB_{AB} + AA_B$	$BB_A$	0	2	2	2	2	0	2	0	0	0	$P_{BA} \rightarrow Q2, P1, T4;$
<i>S</i>													
1	$AA_A + BB_A$		2	2	0	0	0	0	0	0	0	0	
2	$AA_B + BB_B$		0	0	0	2	2	0	0	0	0	0	
3	$AB_{AB} + AB_{BA}$		1	1	2	1	1	0	0	0	0	0	
4	$AA_{AB} + BB_A$	$AB_{BA}$	1	2	1	1	0	0	1	1	1	0	$P_{BA} \rightarrow S1, T5, T1;$
													$T5, T3, P1;$

# Problems with the BWS- and Compeau-Conceptualizations

- ▶ Unwieldy formulas:
  - ▶ Compeau:

$$d_{DCJ}^{ind}(\Pi, \Gamma) = N - [(c + p^{\pi, \pi} + p^{\gamma, \gamma} + \lfloor \frac{p^{\pi, \gamma}}{2} \rfloor) + \frac{1}{2}(p_{even}^0 + \min\{p_{odd}^{\pi}, p_{even}^{\pi}\} + \min\{p_{odd}^{\gamma}, p_{even}^{\gamma}\} + \delta)]$$



# Problems with the BWS- and Compeau-Conceptualizations

- Unwieldy formulas:
  - Compeau:

$$d_{DCJ}^{ind}(\Pi, \Gamma) = N - [(c + p^{\pi, \pi} + p^{\gamma, \gamma} + \lfloor \frac{p^{\pi, \gamma}}{2} \rfloor) + \frac{1}{2}(p_{even}^0 + \min\{p_{odd}^{\pi}, p_{even}^{\pi}\} + \min\{p_{odd}^{\gamma}, p_{even}^{\gamma}\} + \delta)]$$

$\delta = 1$  when  $p^{\pi, \gamma}$  odd and either  $p_{odd}^{\pi} > p_{even}^{\pi}, p_{odd}^{\gamma} > p_{even}^{\gamma}$  or  $p_{odd}^{\pi} < p_{even}^{\pi}, p_{odd}^{\gamma} < p_{even}^{\gamma}$ ; otherwise  $\delta = 0$ .

# Problems with the BWS- and Compeau-Conceptualizations

- Unwieldy formulas:

- BWS:

$$d_{DCJ}^{id}(A, B) = d_{DCJ}(A, B) + \sum_C \lambda(C) \\ - 2P - 3Q - 2T - S - 2M - N$$

- Compeau:

$$d_{DCJ}^{ind}(\Pi, \Gamma) = N - [(c + p^{\pi, \pi} + p^{\gamma, \gamma} + \lfloor \frac{p^{\pi, \gamma}}{2} \rfloor) \\ + \frac{1}{2}(p_{even}^0 + \min\{p_{odd}^{\pi}, p_{even}^{\pi}\} + \min\{p_{odd}^{\gamma}, p_{even}^{\gamma}\} + \delta)]$$

# Problems with the BWS- and Compeau-Conceptualizations

- ▶ Unwieldy formulas:

- ▶ BWS:

$$d_{DCJ}^{id}(A, B) = d_{DCJ}(A, B) + \sum_C \lambda(C) \\ - 2P - 3Q - 2T - S - 2M - N$$

- ▶ Compeau:

$$d_{DCJ}^{ind}(\Pi, \Gamma) = N - [(c + p^{\pi, \pi} + p^{\gamma, \gamma} + \lfloor \frac{p^{\pi, \gamma}}{2} \rfloor) \\ + \frac{1}{2}(p_{even}^0 + \min\{p_{odd}^{\pi}, p_{even}^{\pi}\} + \min\{p_{odd}^{\gamma}, p_{even}^{\gamma}\} + \delta)]$$

- ▶ Models known to be equivalent, but unclear how formulas relate.

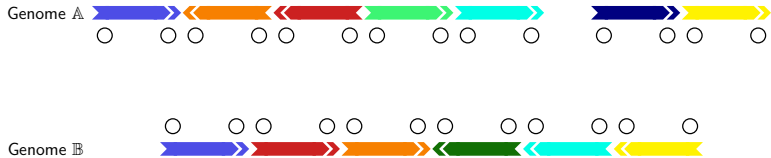
# The Relational Diagram



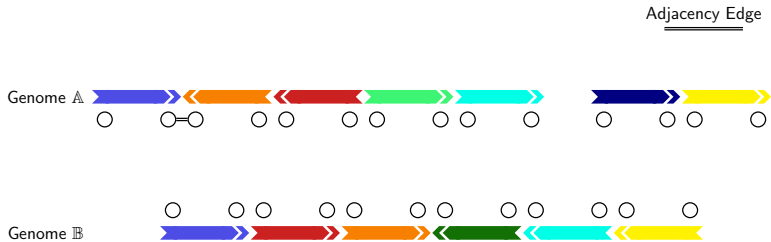
# The Relational Diagram



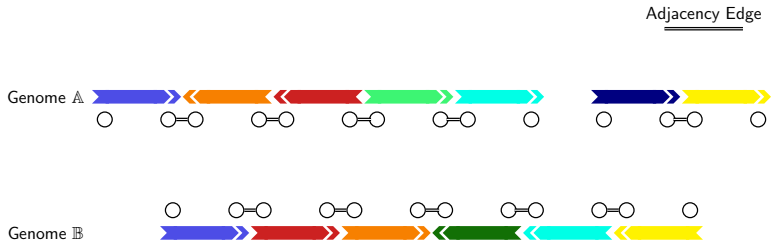
# The Relational Diagram



# The Relational Diagram

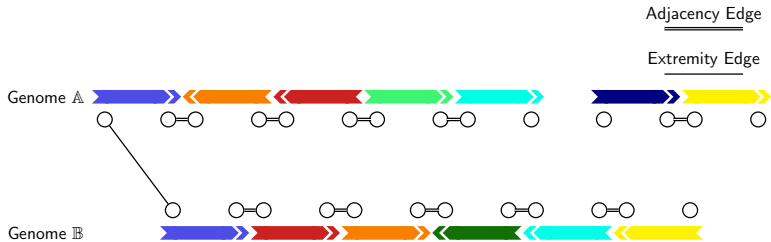


# The Relational Diagram

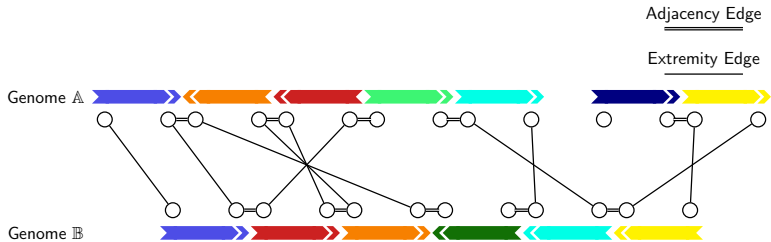




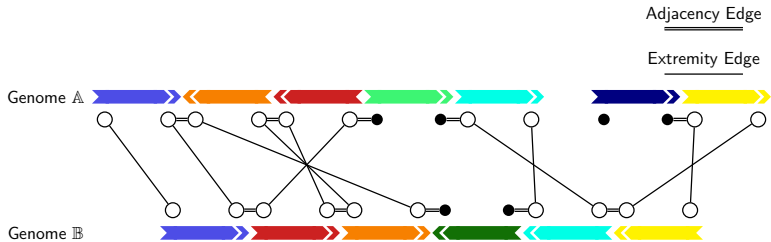
# The Relational Diagram



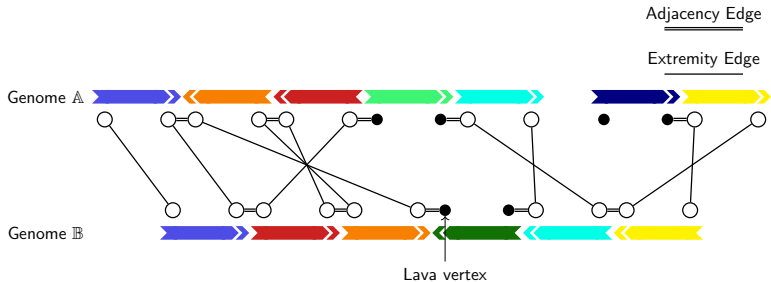
# The Relational Diagram



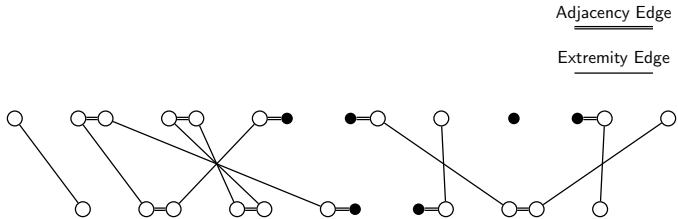
# The Relational Diagram



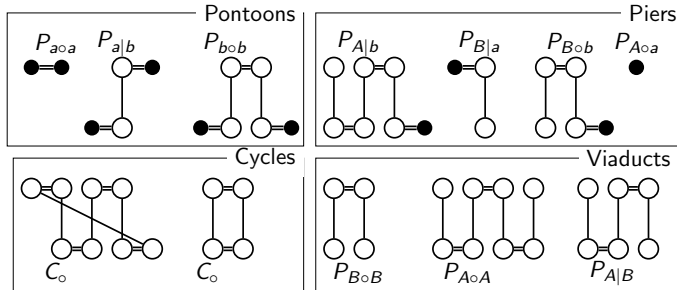
# The Relational Diagram



# The Relational Diagram

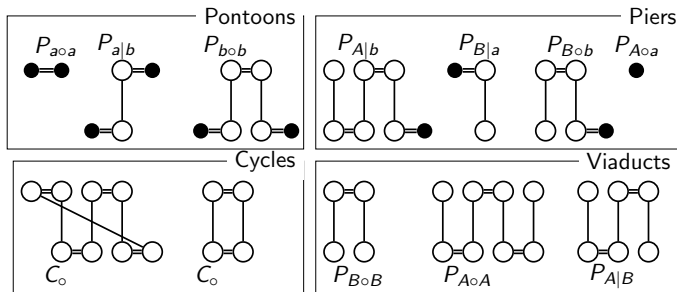


# Types of Components



Earlier work, Recomb-CG 2023.

# A New Formula



Earlier work, Recomb-CG 2023.

$$d_{DCJ}^{id}(\mathbb{A}, \mathbb{B}) = n - c_{\circ} + \left\lceil \frac{p_{a|b} + \max(p_{A \circ a}, p_{B|a}) + \max(p_{A|b}, p_{B \circ b}) - p_{A|B}}{2} \right\rceil$$

# Safe Operations

$$d_{DCJ}^{id}(\mathbb{A}, \mathbb{B}) = n - c_o + \left\lceil \frac{p_{a|b} + \max(p_{A \circ a}, p_{B|a}) + \max(p_{A|b}, p_{B \circ b}) - p_{A|B}}{2} \right\rceil$$



# Safe Operations

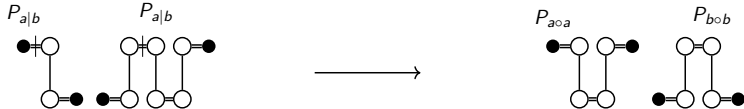
$$d_{DCJ}^{id}(\mathbb{A}, \mathbb{B}) = n - c_o + \left\lceil \frac{p_{a|b} + \max(p_{A \circ a}, p_{B|a}) + \max(p_{A|b}, p_{B \circ b}) - p_{A|B}}{2} \right\rceil$$

Safe operations reduce the formula by 1, no matter the rest of the graph.

# Safe Operations

$$d_{DCJ}^{id}(\mathbb{A}, \mathbb{B}) = n - c_o + \left\lceil \frac{p_{a|b} + \max(p_{A \circ a}, p_{B|a}) + \max(p_{A|b}, p_{B \circ b}) - p_{A|B}}{2} \right\rceil$$

Safe operations reduce the formula by 1, no matter the rest of the graph.

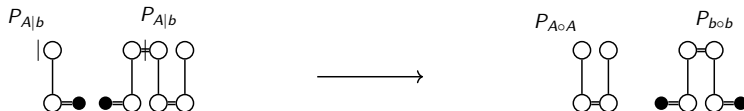


Safe operation

# Safe Operations

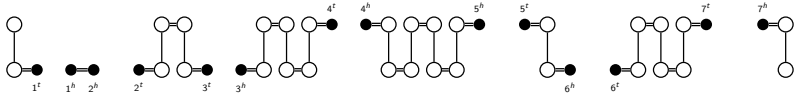
$$d_{DCJ}^{id}(\mathbb{A}, \mathbb{B}) = n - c_o + \left\lceil \frac{p_{a|b} + \max(p_{A \circ a}, p_{B|a}) + \max(p_{A|b}, p_{B \circ b}) - p_{A|B}}{2} \right\rceil$$

Safe operations reduce the formula by 1, no matter the rest of the graph.



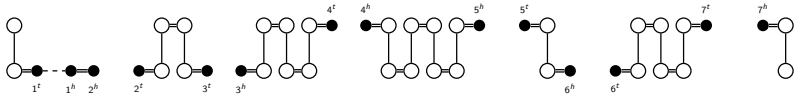
Unsafe operation

# BWS-Conceptualization: Basics

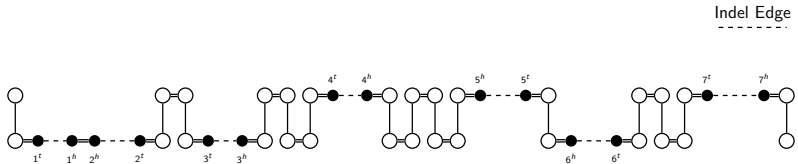


# BWS-Conceptualization: Basics

Indel Edge  
-----

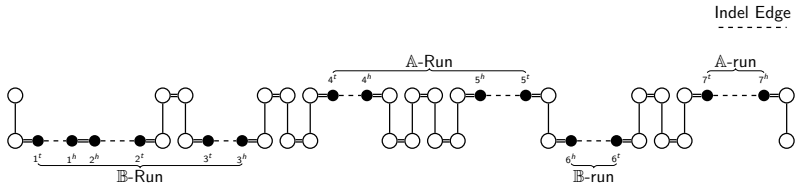


# BWS-Conceptualization: Basics



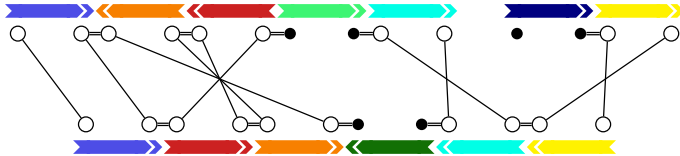


# BWS-Conceptualization: Basics

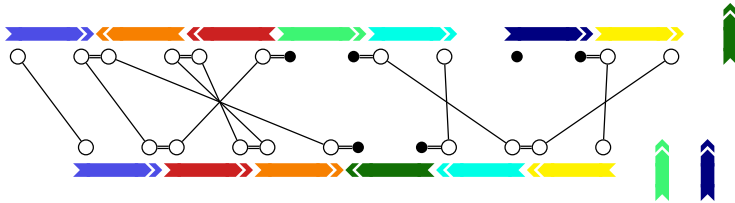




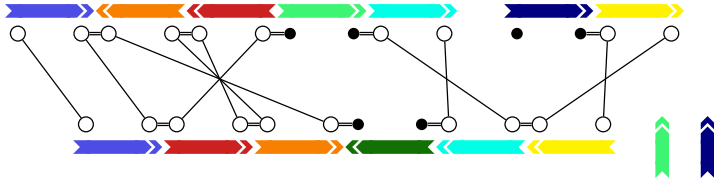
# Compeau-Conceptualization: Basics



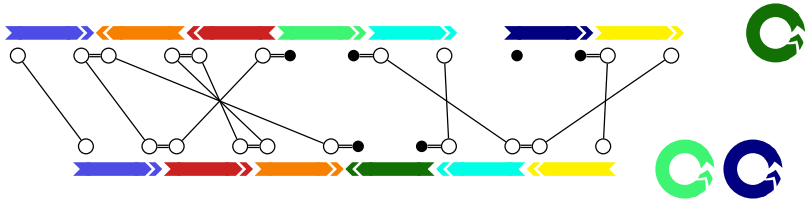
# Compeau-Conceptualization: Basics



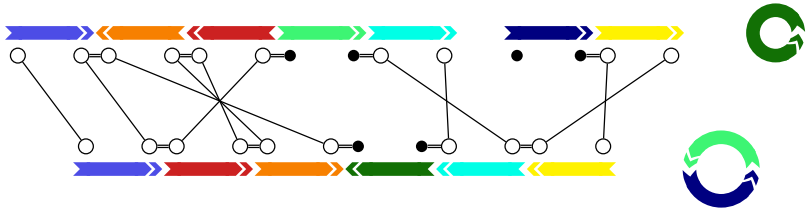
# Compeau-Conceptualization: Basics



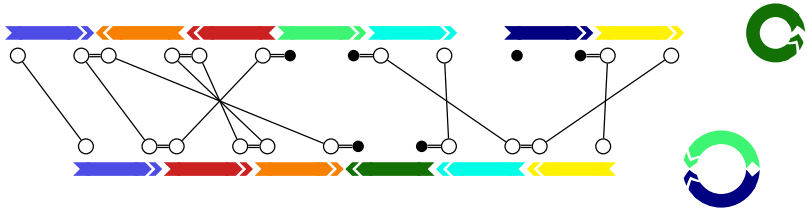
# Compeau-Conceptualization: Basics



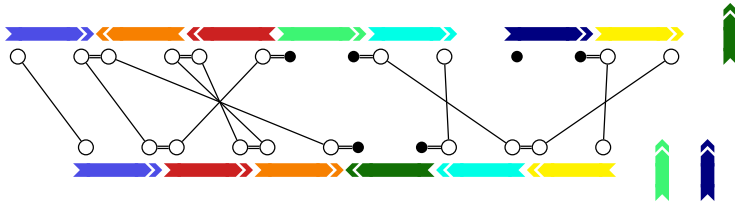
# Compeau-Conceptualization: Basics



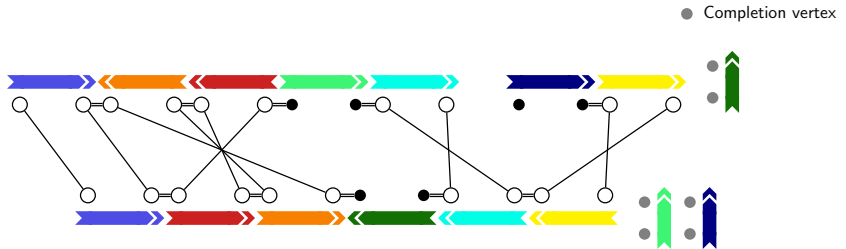
# Compeau-Conceptualization: Basics



# Compeau-Conceptualization: Basics

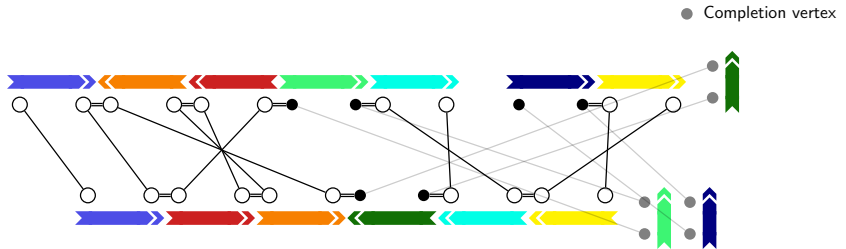


# Compeau-Conceptualization: Basics

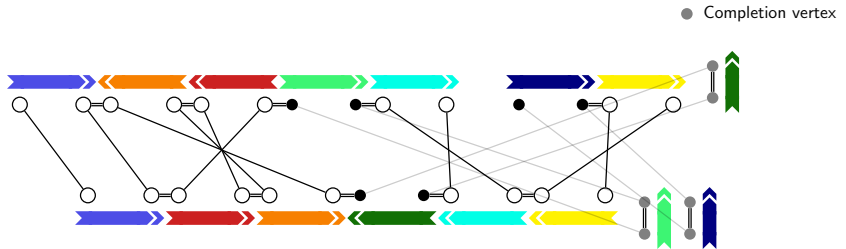




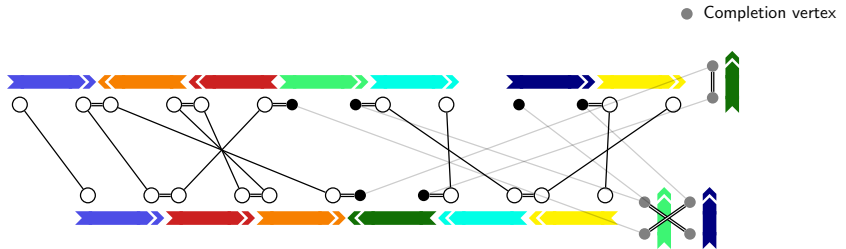
# Compeau-Conceptualization: Basics



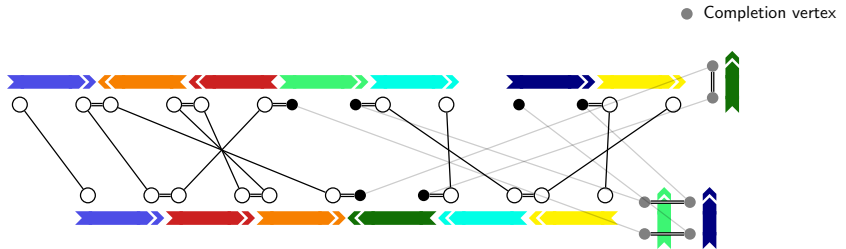
# Compeau-Conceptualization: Basics



# Compeau-Conceptualization: Basics

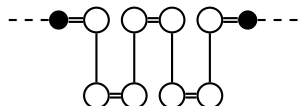


# Compeau-Conceptualization: Basics



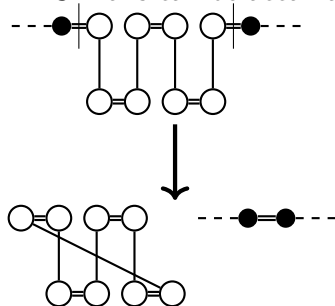
# Bridging BWS- and Compeau-Conceptualizations

BWS: Runs can be accumulated.



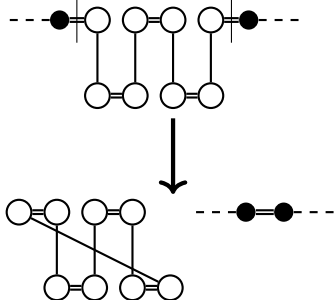
# Bridging BWS- and Compeau-Conceptualizations

BWS: Runs can be accumulated.

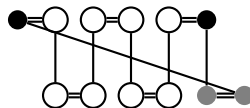


# Bridging BWS- and Compeau-Conceptualizations

BWS: Runs can be accumulated.

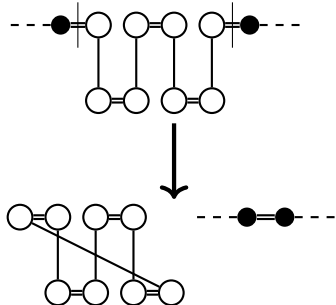


Compeau: Optimal 1-bracelet

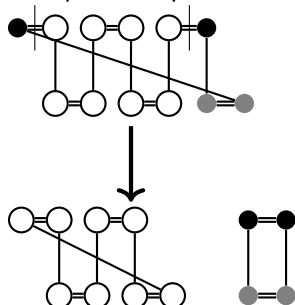


# Bridging BWS- and Compeau-Conceptualizations

BWS: Runs can be accumulated.



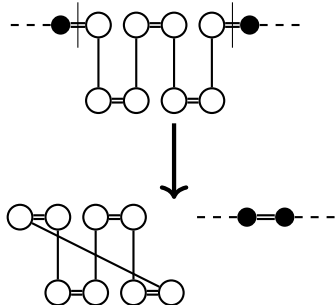
Compeau: Optimal 1-bracelet



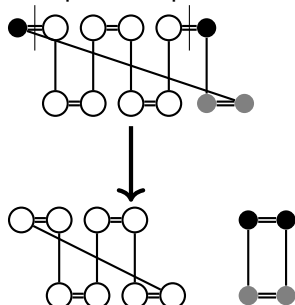


# Bridging BWS- and Compeau-Conceptualizations

BWS: Runs can be accumulated.



Compeau: Optimal 1-bracelet



Common Ground: Safe operation!

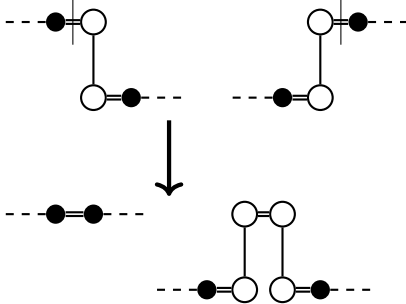
# Bridging BWS- and Compeau-Conceptualizations

BWS: Runs can be recombined.



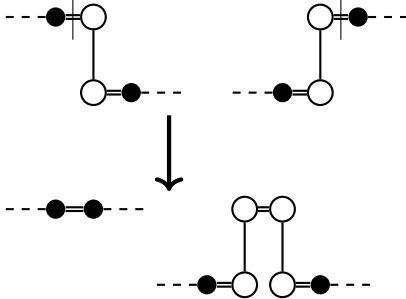
# Bridging BWS- and Compeau-Conceptualizations

BWS: Runs can be recombined.

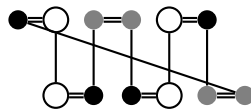


# Bridging BWS- and Compeau-Conceptualizations

BWS: Runs can be recombined.

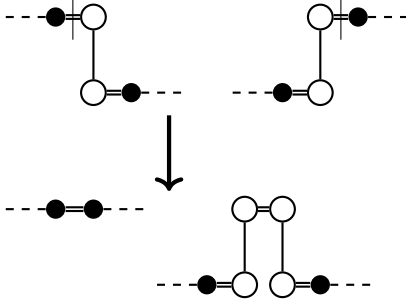


Compeau: Optimal 2-bracelet.

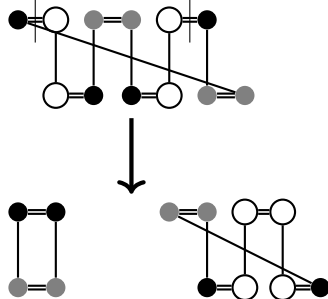


# Bridging BWS- and Compeau-Conceptualizations

BWS: Runs can be recombined.

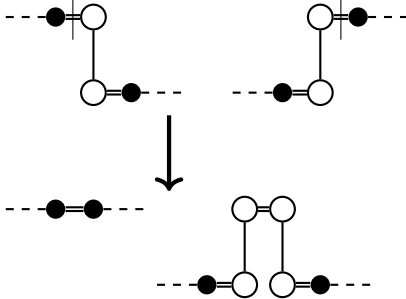


Compeau: Optimal 2-bracelet.

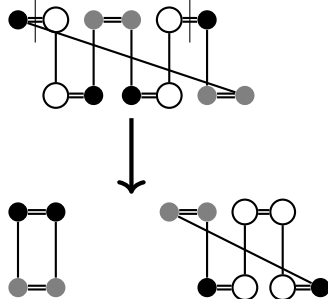


# Bridging BWS- and Compeau-Conceptualizations

BWS: Runs can be recombined.



Compeau: Optimal 2-bracelet.

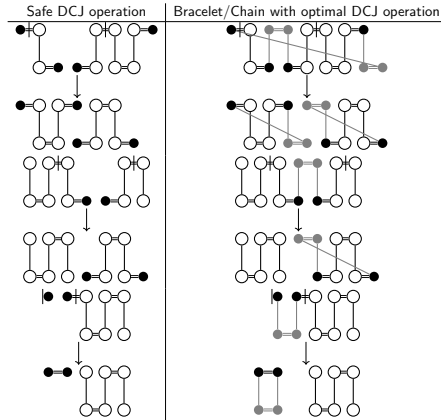


Common Ground: Safe operation!

# Safe Operations and the BWS-Conceptualization

	Sources	Reusable Resultants	Unsat. after safe internal ops.					Unsat. after all safe ops. in group					No partner due to
			$A \circ a$	$B a$	$a b$	$A b$	$B \circ b$	$A \circ a$	$B a$	$P_{a b}$	$A b$	$B \circ b$	
<i>P</i>													
1	$AA_{AB} + BB_{AB}$		1	1	2	1	1	0	0	0	0	0	
<i>Q</i>													
1	$2AA_{AB} + BB_A + BB_B$		2	2	2	2	2	0	0	0	0	0	
2	$2BB_{AB} + AA_A + AA_B$		2	2	2	2	2	0	0	0	0	0	
<i>T</i>													
1	$AA_{AB} + BB_A + AB_{AB}$		2	2	2	1	1	0	0	0	0	0	
2	$AA_{AB} + BB_B + AB_{BA}$		1	1	2	2	2	0	0	0	0	0	
3	$BB_{AB} + AA_A + AB_{BA}$		2	2	2	1	1	0	0	0	0	0	
4	$BB_{AB} + AA_B + AB_{AB}$		1	1	2	2	2	0	0	0	0	0	
5	$2AA_{AB} + BB_A$	$AA_B$	2	2	2	2	0	0	0	0	2	0	$P_{A b} \rightarrow T1, Q1, P1;$ $P_{A \circ a} \rightarrow$ $T2, Q1, P1;$ $P_{B \circ b} \rightarrow$ $Q2, P1, T3;$ $P_{B a} \rightarrow Q2, P1, T4;$
6	$2AA_{AB} + BB_B$	$AA_A$	2	0	2	2	2	2	0	0	0	0	
7	$2BB_{AB} + AA_A$	$BB_B$	2	2	2	0	2	0	0	0	0	2	
8	$2BB_{AB} + AA_B$	$BB_A$	0	2	2	2	2	0	2	0	0	0	
<i>S</i>													
1	$AA_A + BB_A$		2	2	0	0	0	0	0	0	0	0	
2	$AA_B + BB_B$		0	0	0	2	2	0	0	0	0	0	
3	$AB_{AB} + AB_{BA}$		1	1	2	1	1	0	0	0	0	0	
4	$AA_{AB} + BB_A$	$AB_{BA}$	1	2	1	1	0	0	1	1	1	0	$P_{B a} \rightarrow S1, T5, T1;$ $*P_{a b} \rightarrow$ $T5, T1, P1;$

# Safe Operations and the Compeau-Conceptualization

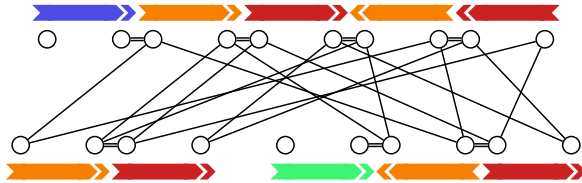




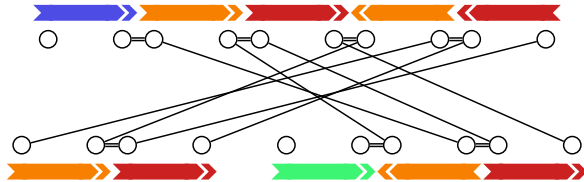
# What's the use?

(Besides didactics)

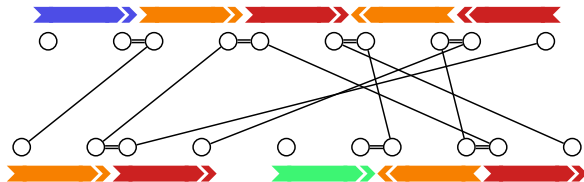
# Duplicate Markers make the Problem NP-hard



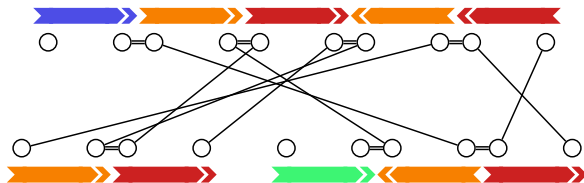
# Duplicate Markers make the Problem NP-hard



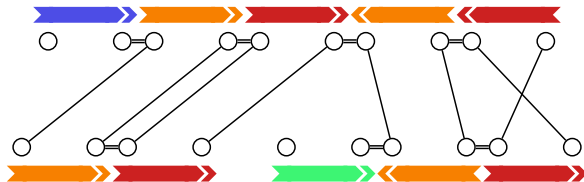
# Duplicate Markers make the Problem NP-hard



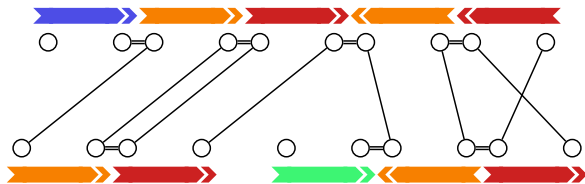
# Duplicate Markers make the Problem NP-hard



# Duplicate Markers make the Problem NP-hard



# Duplicate Markers make the Problem NP-hard



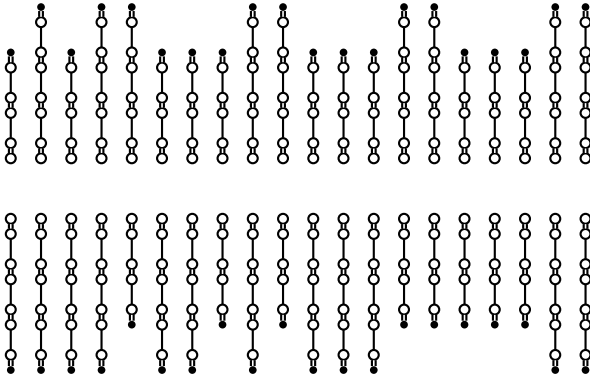
→ Solved by `ding`, ILP based on BWS-model (BBDS,2023).

# Path-Recombinations in BWS-model are too Complicated for an ILP.

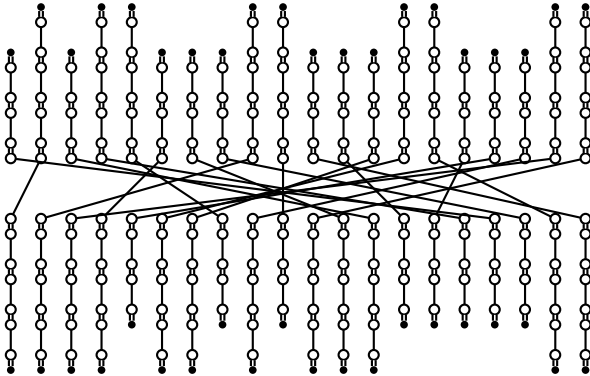
	Sources	Reusable Resultants	Unsat. after safe internal ops.			Unsat. after all safe ops. in group			No partner due to		
			$A \leftrightarrow B$	$B \rightarrow A$	$A \rightarrow B$	$A \leftrightarrow B$	$B \rightarrow A$	$A \rightarrow B$	$A \rightarrow B$	$B \rightarrow A$	
$P$											
1	$AA_{AB} + BB_{AB}$		1	1	2	1	1	0	0	0	
$Q$											
1	$2AA_{AB} + BB_A + BB_B$		2	2	2	2	2	0	0	0	
2	$2BB_{AB} + AA_A + AA_B$		2	2	2	2	2	0	0	0	
$T$											
1	$AA_{AB} + BB_A + AB_{AB}$		2	2	2	1	1	0	0	0	
2	$AA_{AB} + BB_B + AB_{BA}$		1	1	2	2	2	0	0	0	
3	$BB_{AB} + AA_A + AB_{BA}$		2	2	2	1	1	0	0	0	
4	$BB_{AB} + AA_B + AB_{AB}$		1	1	2	2	2	0	0	0	
5	$2AA_{AB} + BB_A$	$AA_B$	2	2	2	2	0	0	0	2	$P_{AB} \rightarrow T1, Q1, P1;$
6	$2AA_{AB} + BB_B$	$AA_A$	2	0	2	2	2	2	0	0	$P_{A \leftrightarrow B} \rightarrow$
7	$2BB_{AB} + AA_A$	$BB_B$	2	2	2	0	2	0	0	0	$T2, Q1, P1;$
8	$2BB_{AB} + AA_B$	$BB_A$	0	2	2	2	2	0	2	0	$P_{B \rightarrow A} \rightarrow$
$S$											
1	$AA_A + BB_A$		2	2	0	0	0	0	0	0	$Q2, P1, T3;$
2	$AA_B + BB_B$		0	0	0	2	2	0	0	0	$P_{B a} \rightarrow Q2, P1, T4;$
3	$AB_{AB} + AB_{BA}$		1	1	2	1	1	0	0	0	
4	$AA_{AB} + BB_A$	$AB_{BA}$	1	2	1	1	0	0	1	1	$P_{B a} \rightarrow S1, T5, T1;$
											$*P_{a b} \rightarrow$
											$T5, T1, P1;$



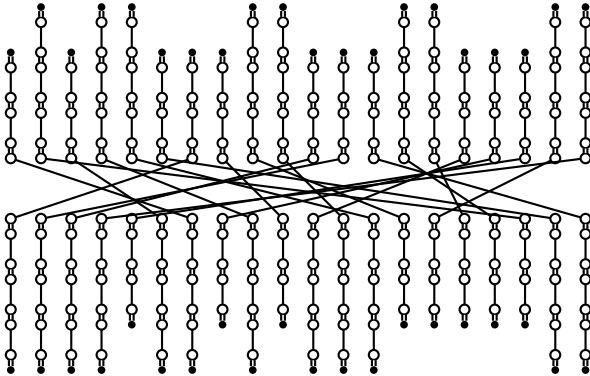
# Transforming Paths into Cycles - Capping



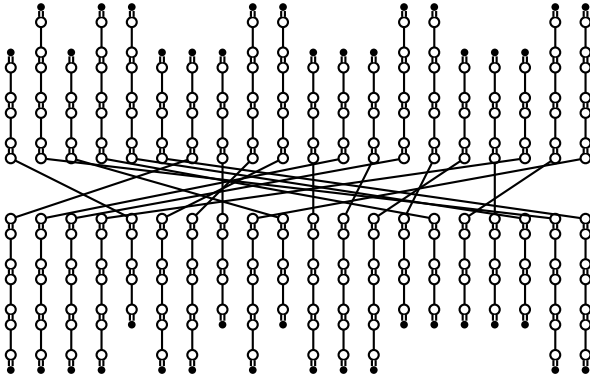
# Transforming Paths into Cycles - Capping



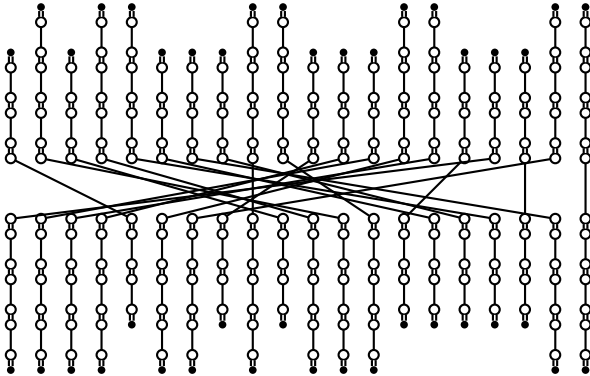
# Transforming Paths into Cycles - Capping



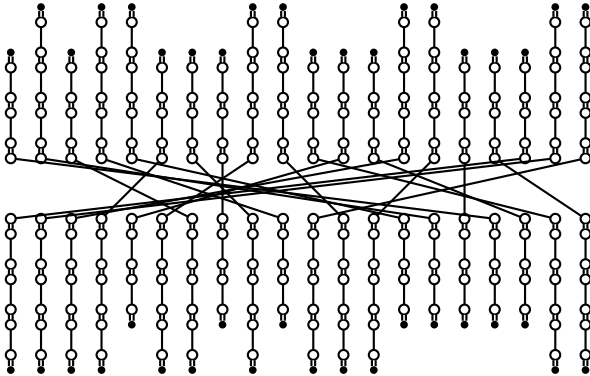
# Transforming Paths into Cycles - Capping



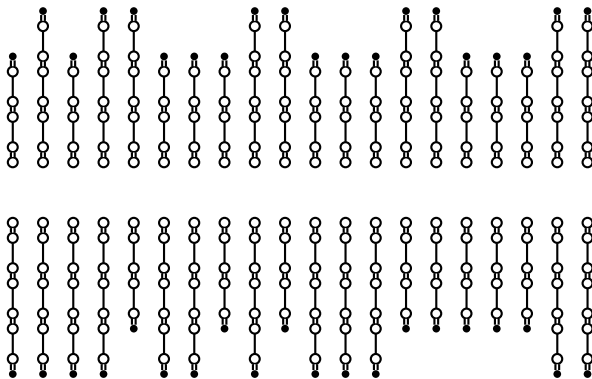
# Transforming Paths into Cycles - Capping



# Transforming Paths into Cycles - Capping



# Transforming Paths into Cycles - Capping



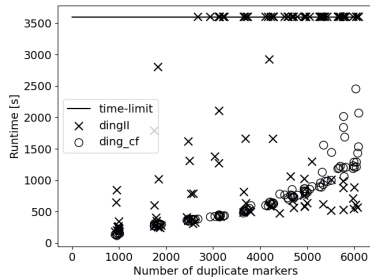
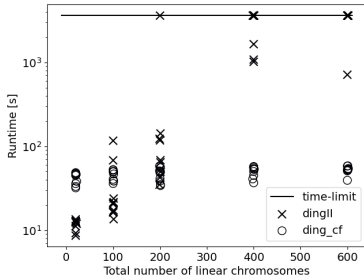
- Superexponential increase of Solution Space. (Rubert & Braga, 2022)

# ILP Based on New Formula Avoids Capping!

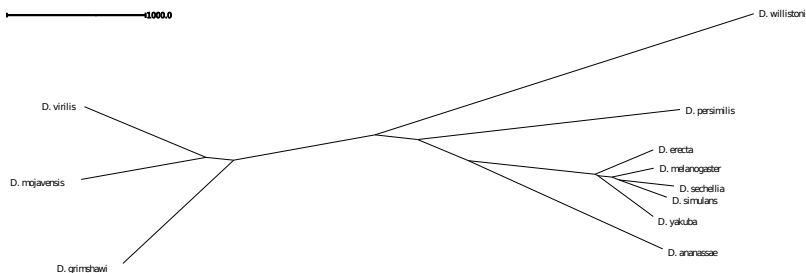
$$d_{DCJ}^{id}(\mathbb{A}, \mathbb{B}) = n - c_{\circ} + \left\lceil \frac{p_{a|b} + \max(p_{A \circ a}, p_{B|a}) + \max(p_{A|b}, p_{B \circ b}) - p_{A|B}}{2} \right\rceil$$



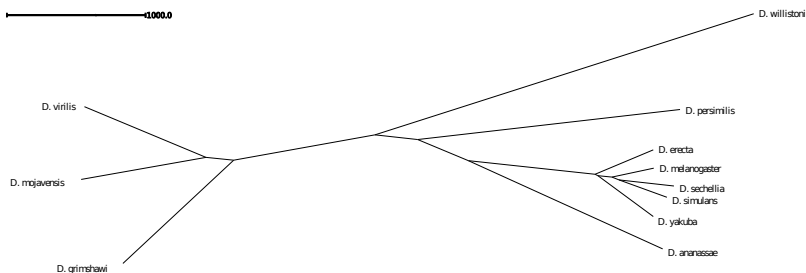
# Considerable Performance Advantage when the Number of Linear Chromosomes is High



# Computing Rearrangement Phylogenies on Contig-level resolved Genomes

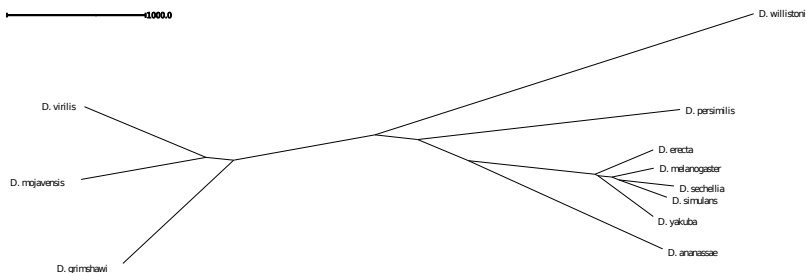


# Computing Rearrangement Phylogenies on Contig-level resolved Genomes



► High quality

# Computing Rearrangement Phylogenies on Contig-level resolved Genomes



- ▶ High quality
- ▶ Robust: On average only 0.53% deviation from computed pairwise distances.

# Summarized Results

- ▶ More compact and simple DCJ-indel formula

# Summarized Results

- ▶ More compact and simple DCJ-indel formula
  - ▶ Link between BWS and Compeau-conceptualizations

# Summarized Results

- ▶ More compact and simple DCJ-indel formula
  - ▶ Link between BWS and Compeau-conceptualizations
- ▶ Significant performance improvement of `ding` for high numbers of linear chromosomes/contigs

