

# The Space of Discrete Coalescent Trees

Lena Collienue

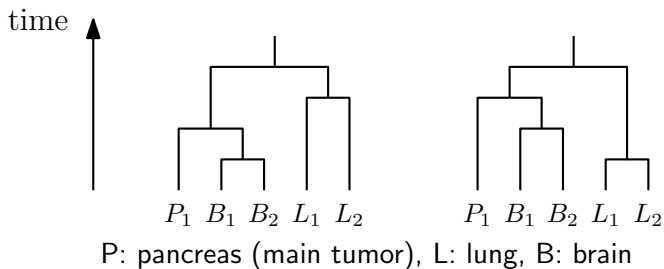


Biological Data Science Lab  
Department of Computer Science  
University of Otago

30/03/2021

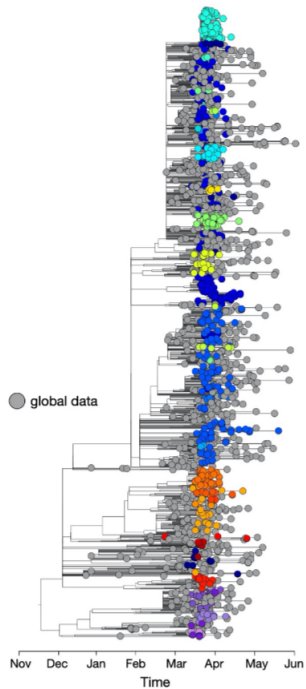
# Time Trees

## Cancer Phylogenies



# Time Trees

## SARS-CoV-2

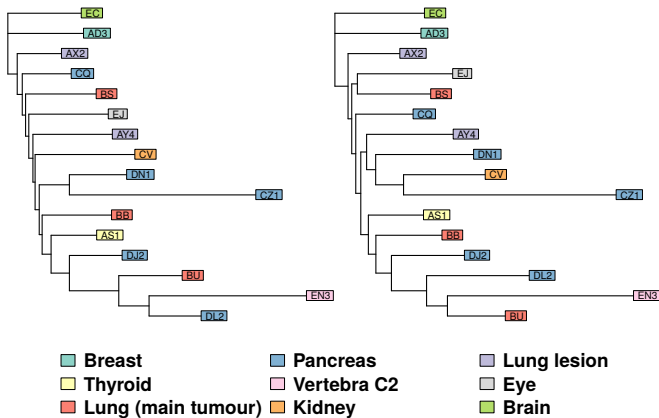


# Why Distances?

Tree Inference

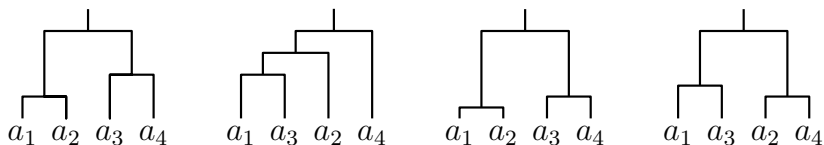
# Why Distances?

## Tree Inference



# Why Distances?

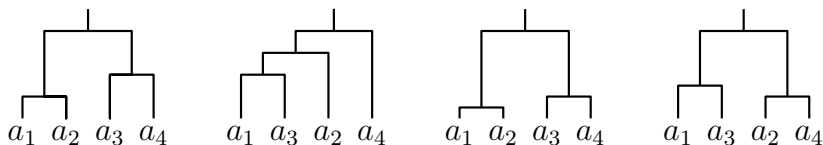
## Summarising Trees



What is the mean tree?

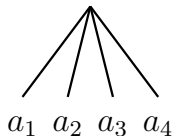
# Why Distances?

## Summarising Trees

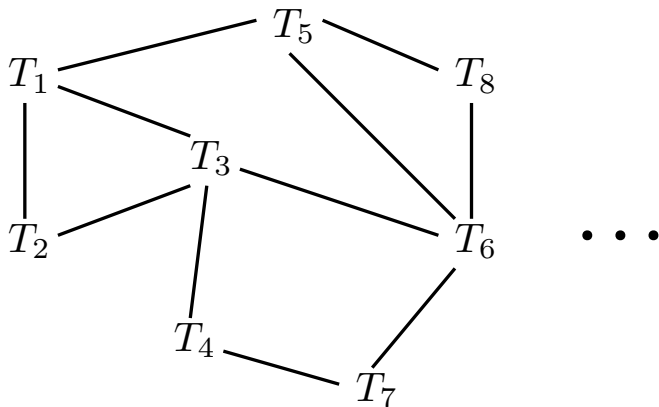


What is the mean tree?

Problem: In most tree spaces the mean tree is a star tree:

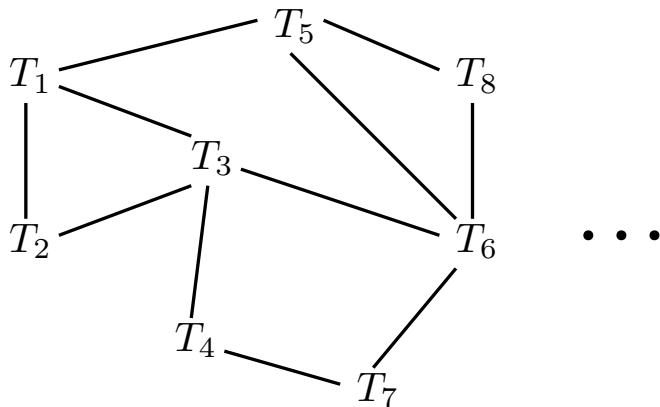


# Discrete Tree Spaces





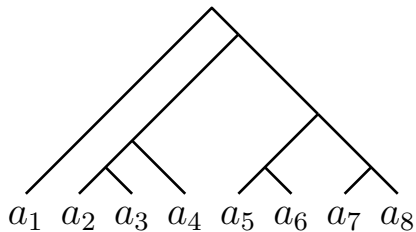
# Discrete Tree Spaces



Tree Rearrangement operations: NNI, SPR, TBR

# Phylogenetic Trees

Rooted, binary

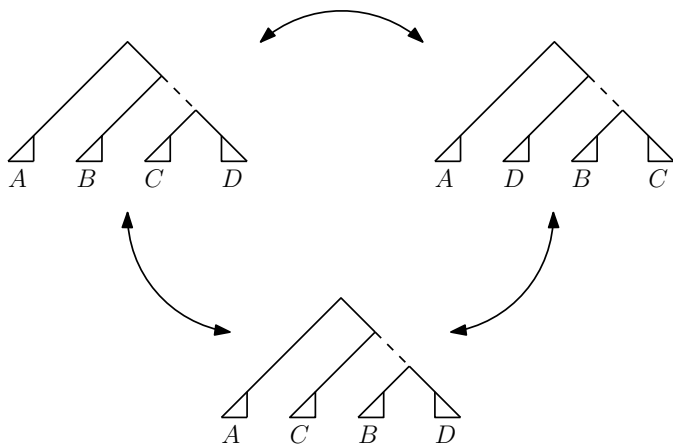


# NNI – Nearest Neighbour Interchange

## Definition 1

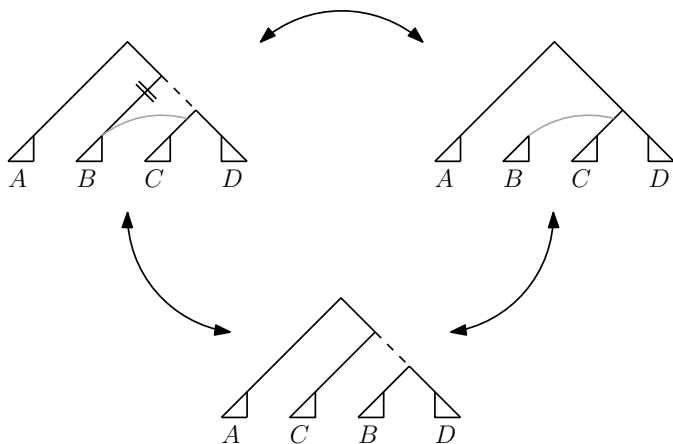
# NNI – Nearest Neighbour Interchange

## Definition 1



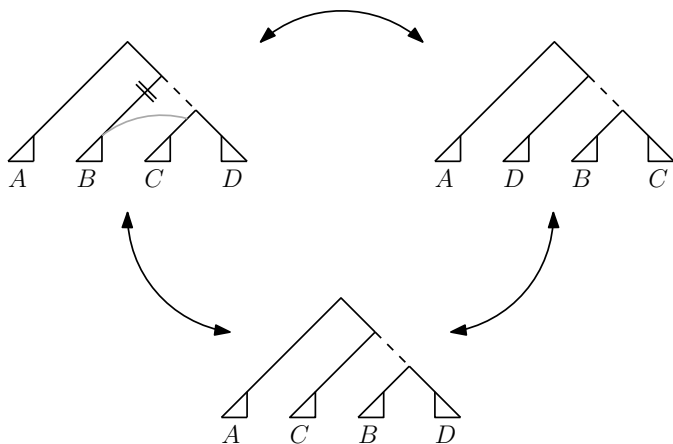
# NNI – Nearest Neighbour Interchange

## Definition 1



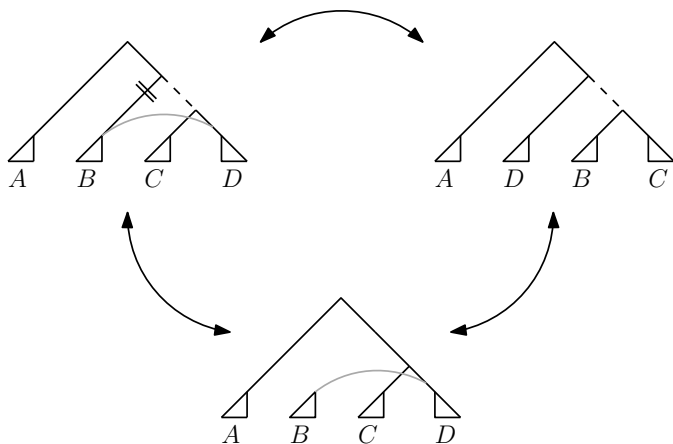
# NNI – Nearest Neighbour Interchange

## Definition 1



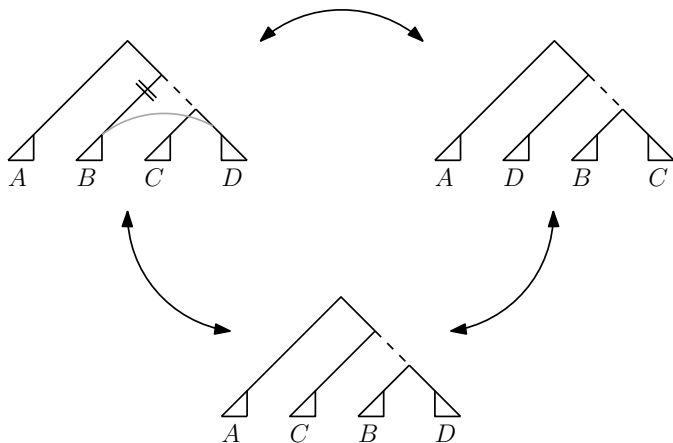
# NNI – Nearest Neighbour Interchange

## Definition 1



# NNI – Nearest Neighbour Interchange

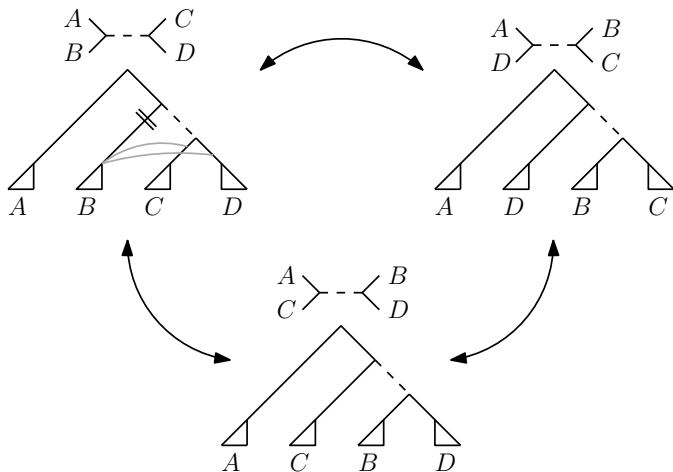
## Definition 1





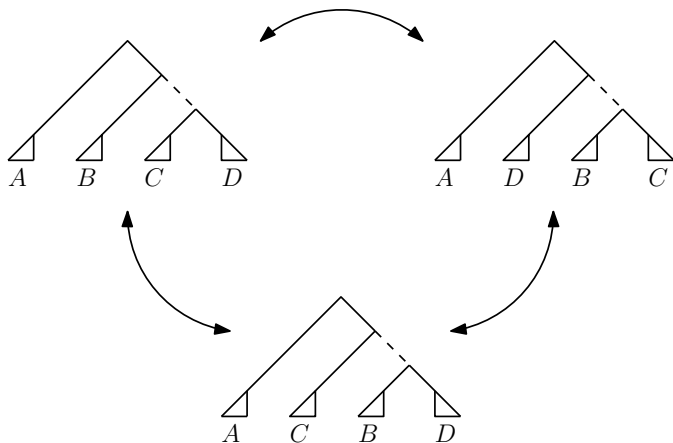
# NNI – Nearest Neighbour Interchange

## Definition 1



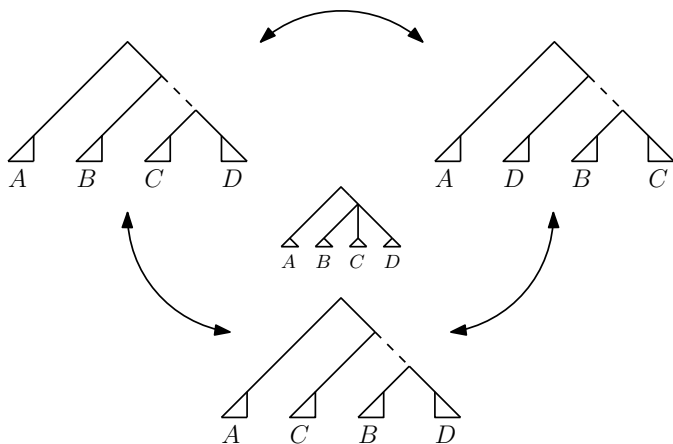
# NNI – Nearest Neighbour Interchange

## Definition 2



# NNI – Nearest Neighbour Interchange

## Definition 2



# Computing Distances

NNI

NNI-DIST:

INSTANCE: A pair of trees  $T$  and  $R$

FIND: Distance between  $T$  and  $R$  in NNI

# Computing Distances

NNI

NNI-DIST:

INSTANCE: A pair of trees  $T$  and  $R$

FIND: Distance between  $T$  and  $R$  in NNI

►  $\mathcal{NP}$ -hard

# Computing Distances

## NNI

### NNI-DIST:

INSTANCE: A pair of trees  $T$  and  $R$

FIND: Distance between  $T$  and  $R$  in NNI

- ▶  $\mathcal{NP}$ -hard
- ▶ BUT: fixed-parameter tractable (FPT):

distance computable in  $\mathcal{O}(2^{\frac{21k}{2}} * n)$  where  $d(T, R) \leq k$

# Computing Distances

## NNI

### NNI-DIST:

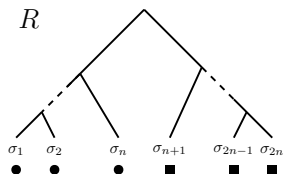
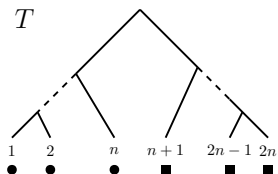
INSTANCE: A pair of trees  $T$  and  $R$

FIND: Distance between  $T$  and  $R$  in NNI

- ▶  $\mathcal{NP}$ -hard
- ▶ BUT: fixed-parameter tractable (FPT):  
distance computable in  $\mathcal{O}(2^{\frac{21k}{2}} * n)$  where  $d(T, R) \leq k$
- ▶ Approximation algorithm: ratio  $\mathcal{O}(\log(n))$

# Biological Interpretability

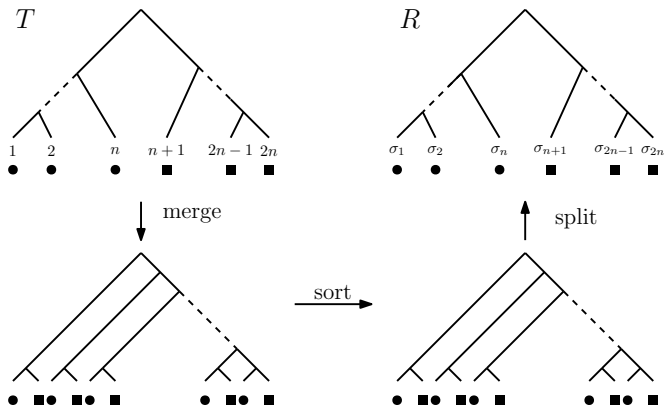
## Cluster Property





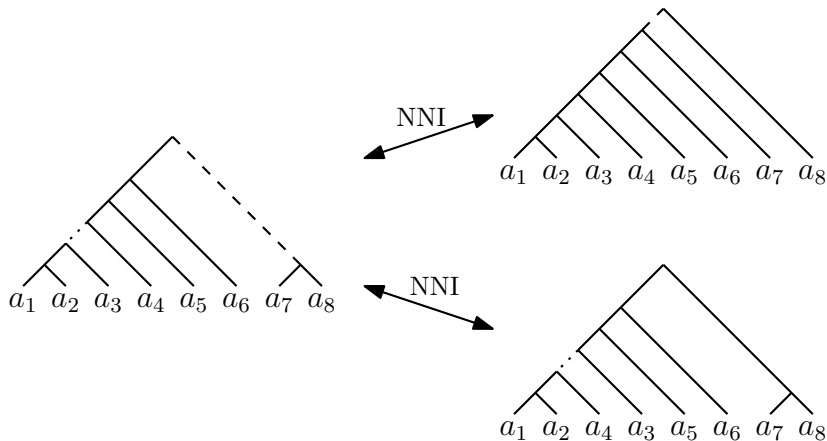
# Biological Interpretability

## Cluster Property

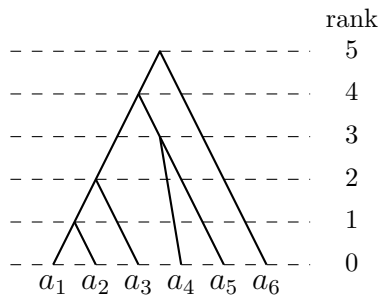


# Biological Interpretability

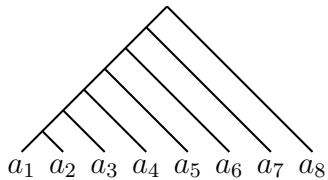
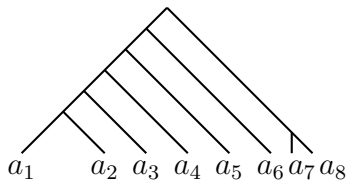
NNI move = NNI move?



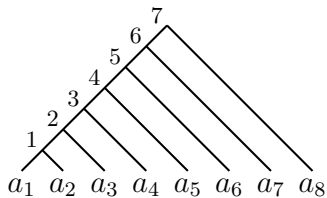
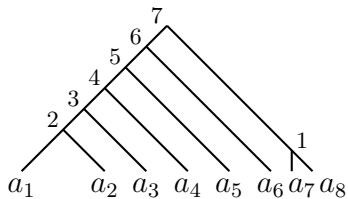
# Ranked Trees



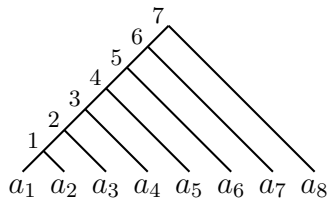
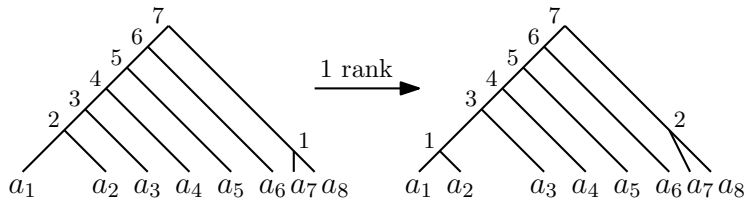
# NNI on Ranked Trees



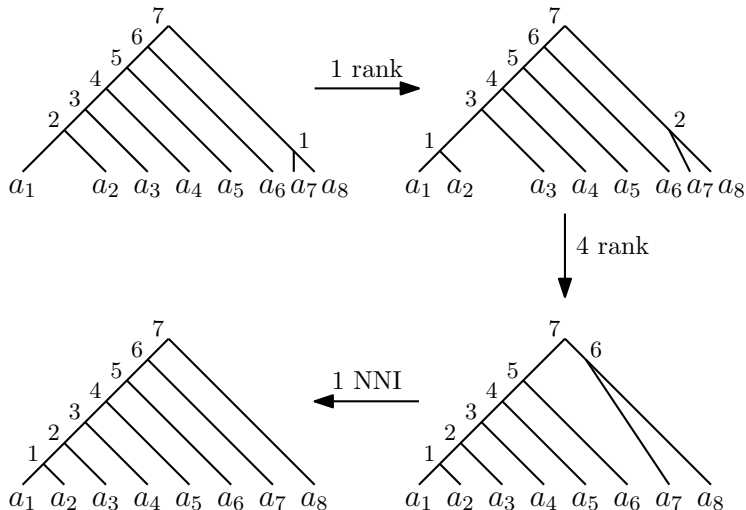
# NNI on Ranked Trees



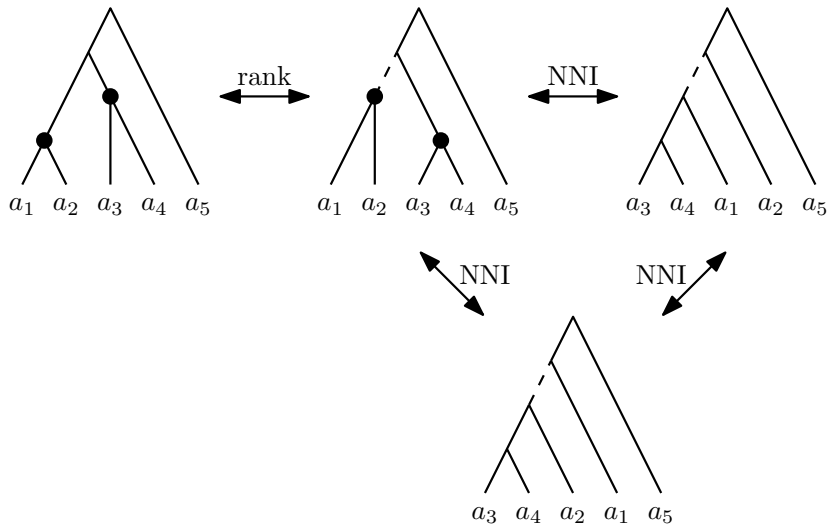
# NNI on Ranked Trees



# NNI on Ranked Trees



# RNNI





# Shortest Path Problem

RNNI

RNNI-SP:

INSTANCE: A pair of ranked trees  $T$  and  $R$

FIND: Shortest Path between  $T$  and  $R$  in RNNI

# FINDPATH

- ▶ Greedy algorithm for approximating RNNI-SP

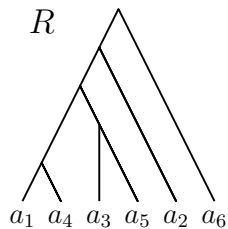
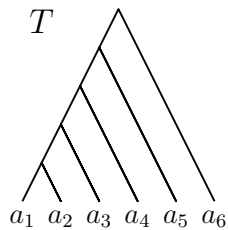
# FINDPATH

- ▶ Greedy algorithm for approximating RNNI-SP
- ▶ Running time  $\mathcal{O}(n^2)$

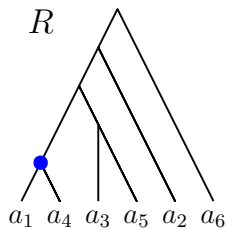
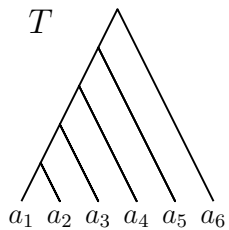
# FINDPATH

- ▶ Greedy algorithm for approximating RNNI-SP
- ▶ Running time  $\mathcal{O}(n^2)$
- ▶ Shortest paths for up to 7 leaves

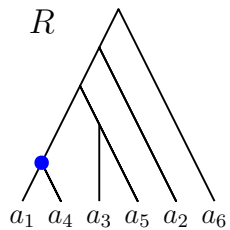
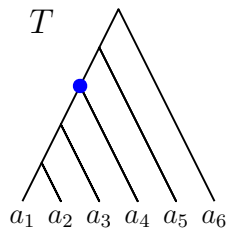
# FINDPATH



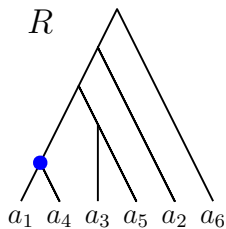
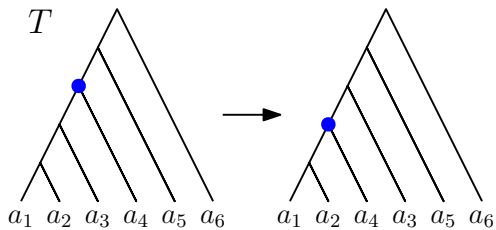
# FINDPATH



# FINDPATH

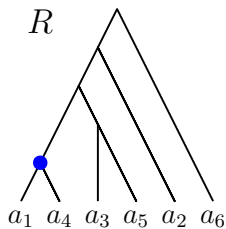
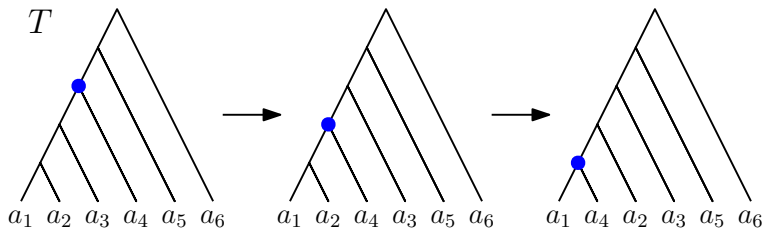


# FINDPATH

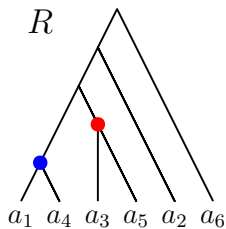
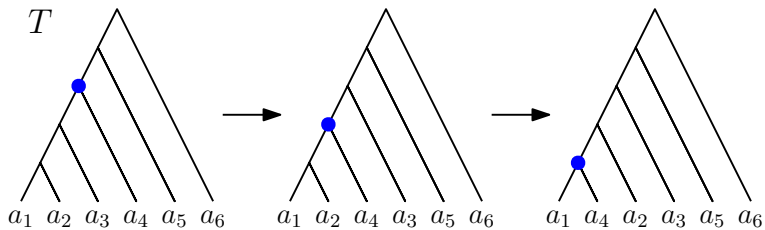




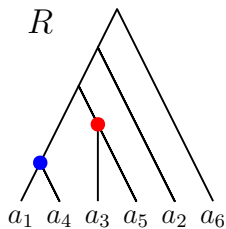
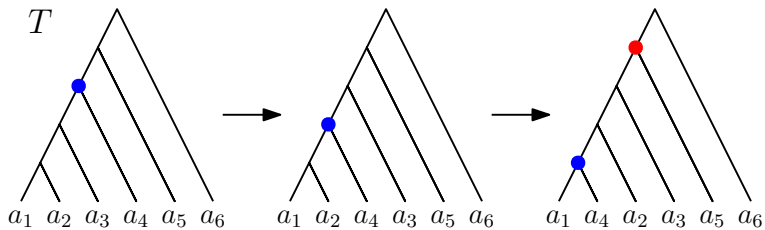
# FINDPATH



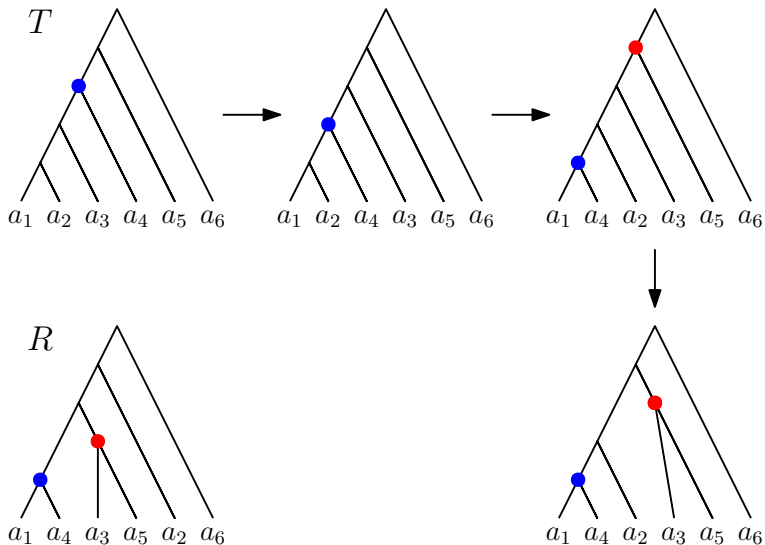
# FINDPATH



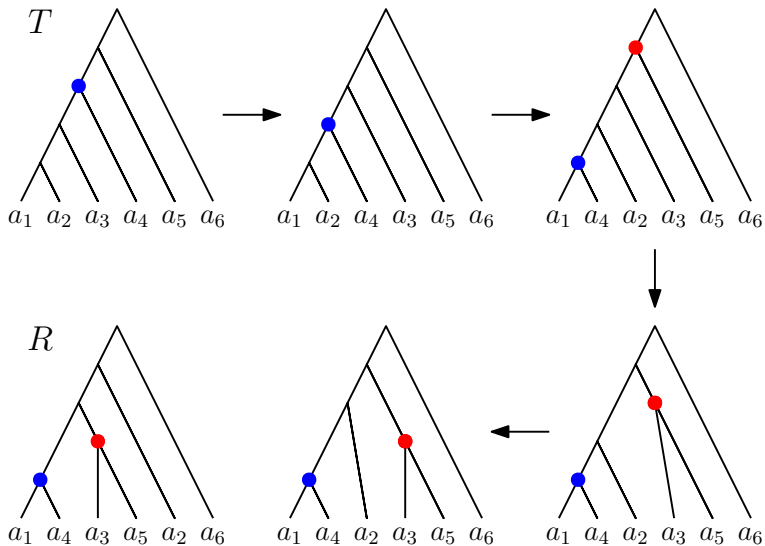
# FINDPATH



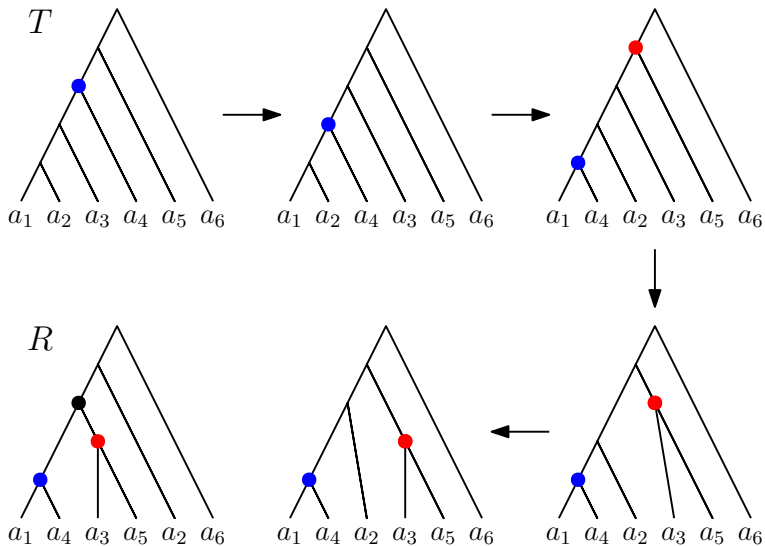
# FINDPATH



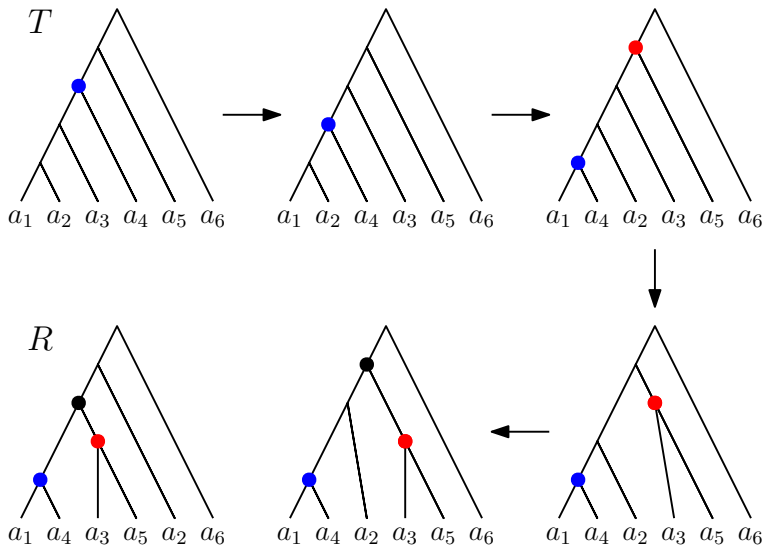
# FINDPATH



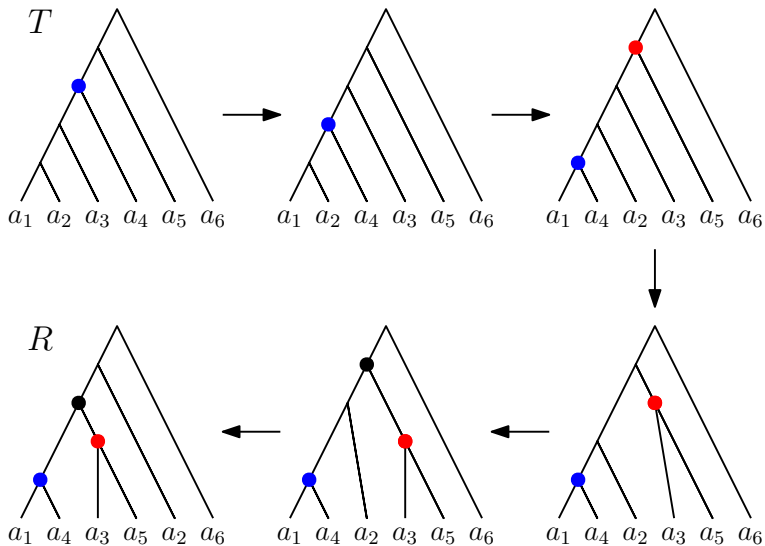
# FINDPATH



# FINDPATH



# FINDPATH





# FINDPATH

## Theorem

FINDPATH *computes shortest paths in* RNNI.

# FINDPATH

## Theorem

FINDPATH *computes shortest paths in* RNNI.

## Idea for proof

$\text{FP}(T, R) :=$  *path between*  $T$  *and*  $R$  *computed by* FINDPATH

# FINDPATH

## Theorem

FINDPATH *computes shortest paths in* RNNI.

## Idea for proof

$\text{FP}(T, R) :=$  *path between*  $T$  *and*  $R$  *computed by* FINDPATH

## Lemma

*If for all trees*  $T$ ,  $R$  *and neighbour*  $T'$  *of*  $T$  *it is*

$$|\text{FP}(T', R)| \geq |\text{FP}(T, R)| - 1,$$

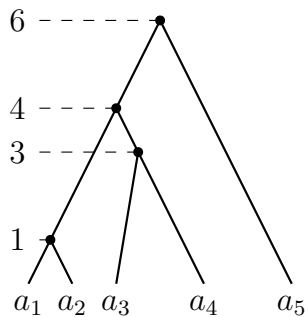
*then*

$$|\text{FP}(T, R)| = d(T, R)$$

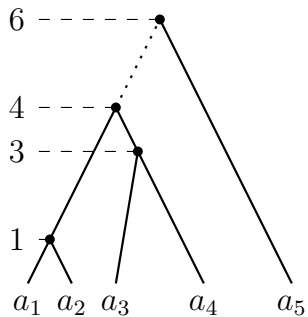
*for all trees*  $T$  *and*  $R$

# Discrete Coalescent Trees

# Discrete Coalescent Trees

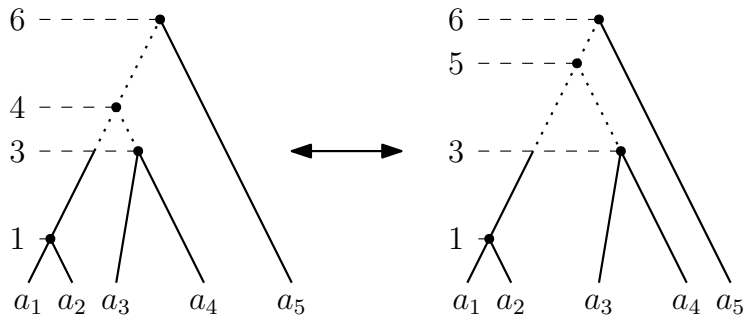


# Discrete Coalescent Trees



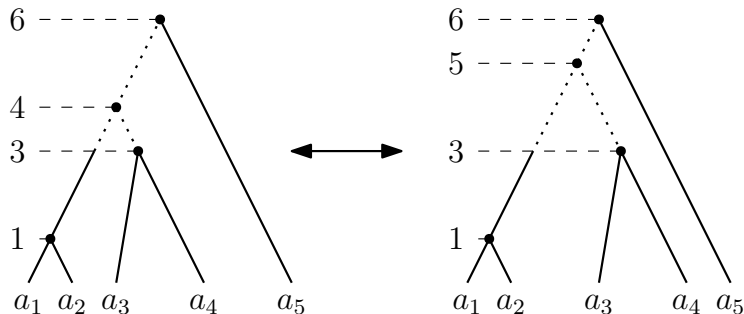
# Discrete Coalescent Trees

Length move



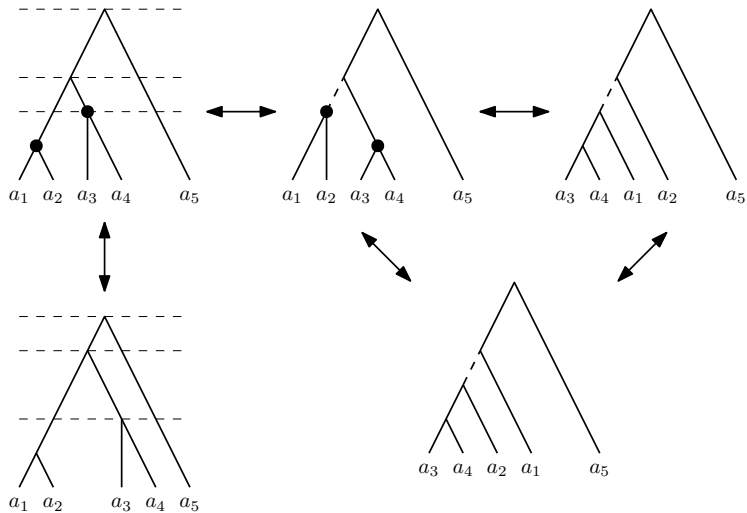
# Discrete Coalescent Trees

Length move

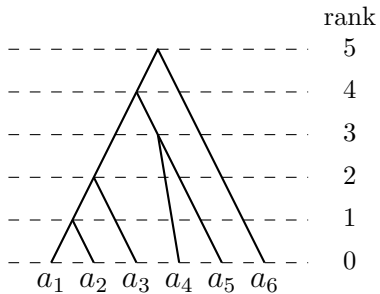


New parameter:  $m = \max$  height of tree





$$\text{DCT}_{n-1} = RNNI$$

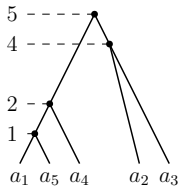
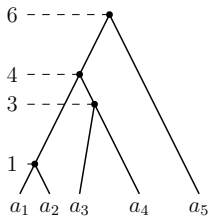


FINDPATH in  $\text{DCT}_m$

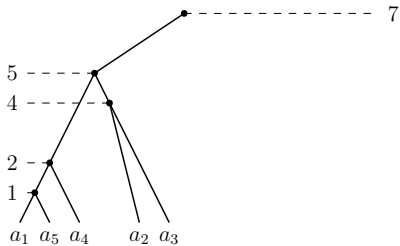
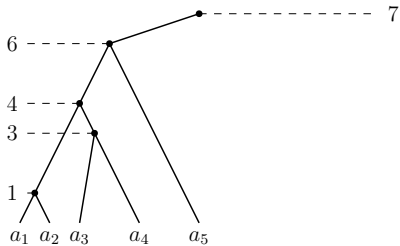
$m = 6$

# FINDPATH in $DCT_m$

$m = 6$

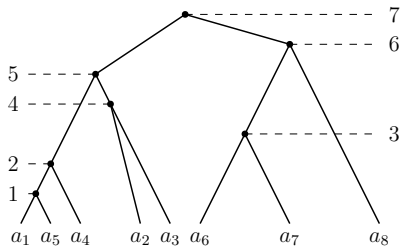
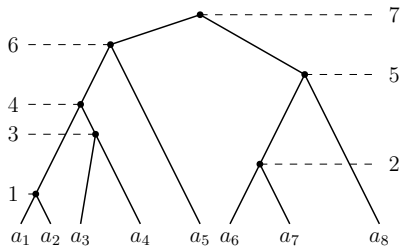


## FINDPATH in $\text{DCT}_m$

$$m = 6$$


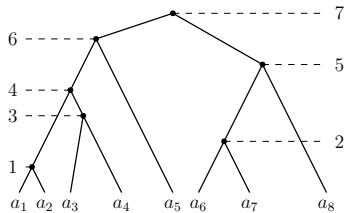
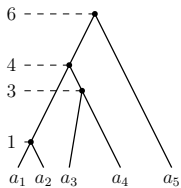
# FINDPATH in $DCT_m$

$m = 6$



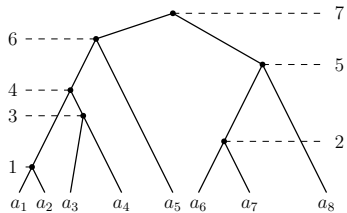
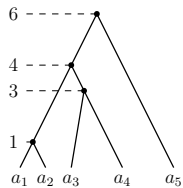
# FINDPATH in $DCT_m$

$m = 6$

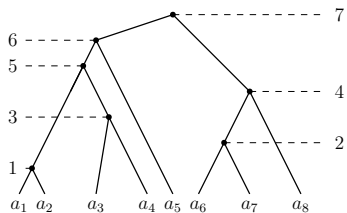


# FINDPATH in $DCT_m$

$m = 6$



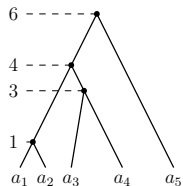
rank move



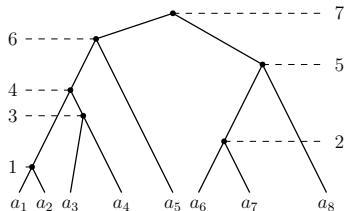
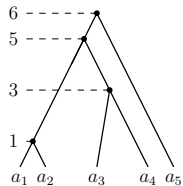


# FINDPATH in $DCT_m$

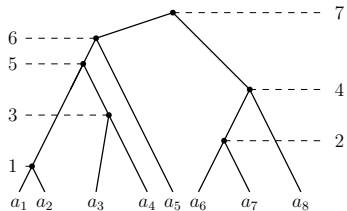
$m = 6$



length move



rank move



# FINDPATH in $\text{DCT}_m$

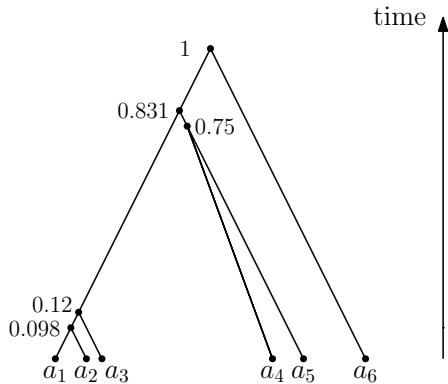
$m = 6$

## Theorem

`FINDPATH` computes shortest paths between discrete coalescent trees  $T$  and  $R$  in  $\mathcal{O}(m^2)$ .

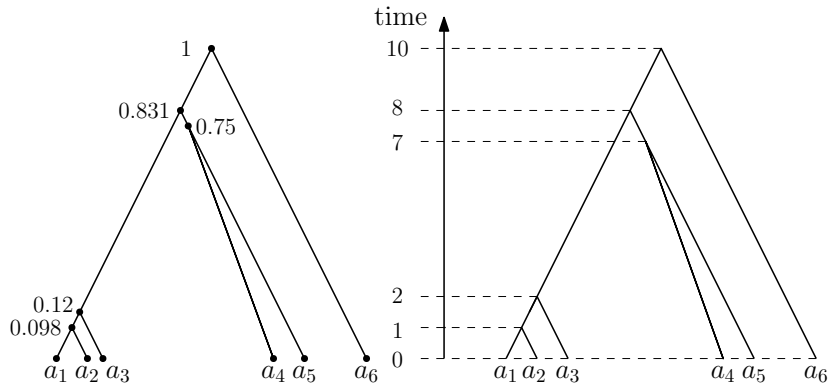
# FINDPATH in $DCT_m$

Scalability



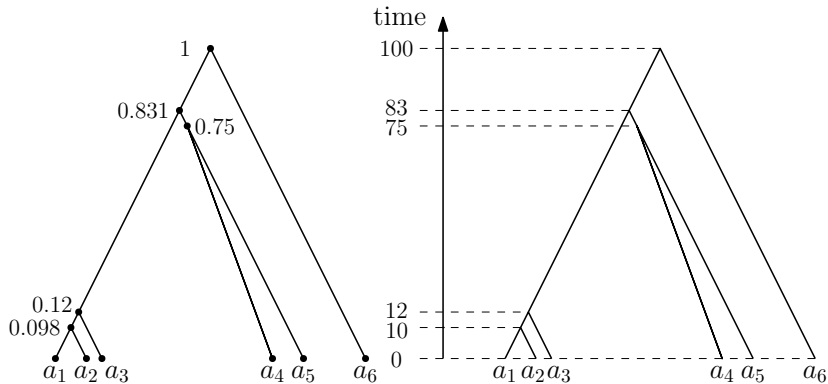
# FINDPATH in $DCT_m$

## Scalability



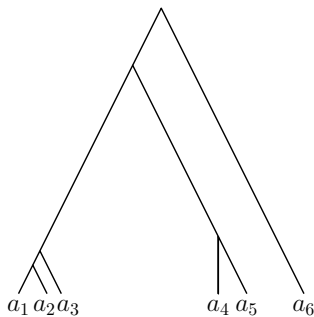
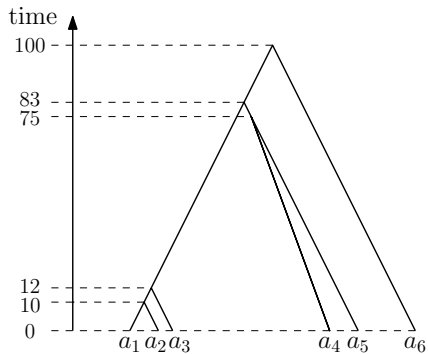
# FINDPATH in $DCT_m$

## Scalability



# FINDPATH in $DCT_m$

## Scalability



## Properties of $\text{DCT}_m$

	# Trees	Diameter	Radius
RNNI	$\frac{n!(n-1)!}{2^{n-1}}$	$\binom{n-1}{2}$	$\binom{n-1}{2}$

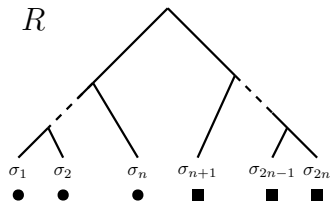
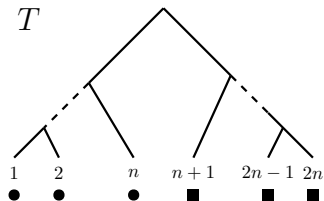
# Properties of $\text{DCT}_m$

	# Trees	Diameter	Radius
RNNI	$\frac{n!(n-1)!}{2^{n-1}}$	$\binom{n-1}{2}$	$\binom{n-1}{2}$
$\text{DCT}_m$	$\frac{n!(n-1)!}{2^{n-1}} \binom{m}{n-1}$	$\binom{n-1}{2} + (m - n + 1)(n - 1)$	?



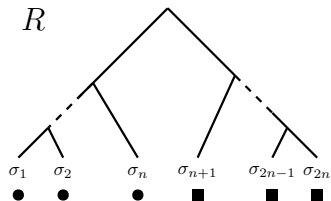
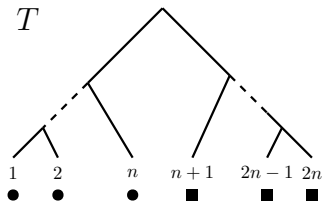
# Properties of $\text{DCT}_m$

## Cluster Property



# Properties of $\text{DCT}_m$

## Cluster Property

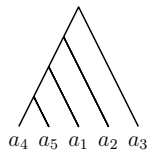
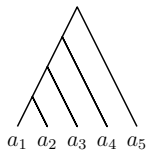


## Theorem

$\text{DCT}_m$  has the cluster property.

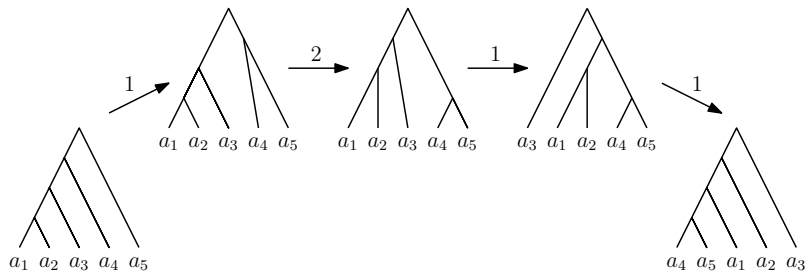
# Properties of $\text{DCT}_m$

## Caterpillar Trees



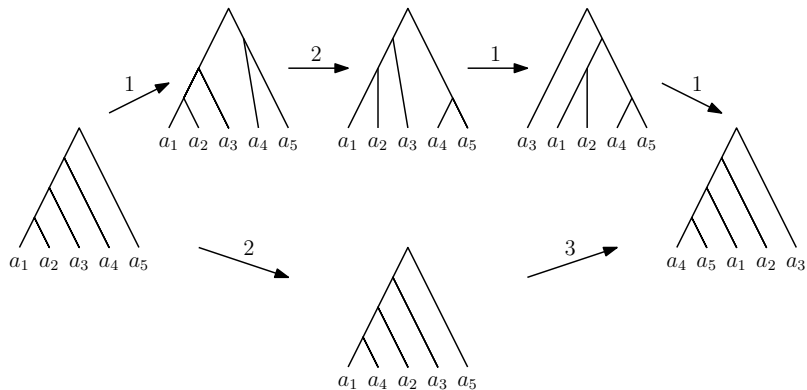
# Properties of $\text{DCT}_m$

## Caterpillar Trees



# Properties of $\text{DCT}_m$

## Caterpillar Trees



# Properties of $\text{DCT}_m$

## Caterpillar Trees

### Theorem

*The set of caterpillar trees is convex.*

# Properties of $\text{DCT}_m$

## Caterpillar Trees

### Theorem

*The set of caterpillar trees is convex.*

### Corollary

*The distance between caterpillar trees can be computed in  $\mathcal{O}(n\sqrt{\log(n)})$ .*

# Conclusion

Solved problems:



# Conclusion

Solved problems:

- ▶ RNNI and  $\text{DCT}_m$ : shortest paths in  $\mathcal{O}(n^2)!$

# Conclusion

Solved problems:

- ▶ RNNI and  $\text{DCT}_m$ : shortest paths in  $\mathcal{O}(n^2)$ !
- ▶ We know diameter, radius, cluster property, convexity of set of caterpillar trees

# Conclusion

Solved problems:

- ▶ RNNI and  $\text{DCT}_m$ : shortest paths in  $\mathcal{O}(n^2)!$
- ▶ We know diameter, radius, cluster property, convexity of set of caterpillar trees

Open Problems:

# Conclusion

Solved problems:

- ▶ RNNI and  $\text{DCT}_m$ : shortest paths in  $\mathcal{O}(n^2)$ !
- ▶ We know diameter, radius, cluster property, convexity of set of caterpillar trees

Open Problems:

- ▶ Can we compute distances more efficiently?

# Conclusion

Solved problems:

- ▶ RNNI and  $\text{DCT}_m$ : shortest paths in  $\mathcal{O}(n^2)$ !
- ▶ We know diameter, radius, cluster property, convexity of set of caterpillar trees

Open Problems:

- ▶ Can we compute distances more efficiently?
- ▶ How can we summarise trees?

# Conclusion

Solved problems:

- ▶ RNNI and  $\text{DCT}_m$ : shortest paths in  $\mathcal{O}(n^2)!$
- ▶ We know diameter, radius, cluster property, convexity of set of caterpillar trees

Open Problems:

- ▶ Can we compute distances more efficiently?
- ▶ How can we summarise trees?
- ▶ Does this help us doing statistics in tree space? Confidence intervals?

# Conclusion

Solved problems:

- ▶ RNNI and  $\text{DCT}_m$ : shortest paths in  $\mathcal{O}(n^2)!$
- ▶ We know diameter, radius, cluster property, convexity of set of caterpillar trees

Open Problems:

- ▶ Can we compute distances more efficiently?
- ▶ How can we summarise trees?
- ▶ Does this help us doing statistics in tree space? Confidence intervals?
- ▶ ...

# Thank you

- ▶ Alex Gavryushkin (University of Otago)
- ▶ David Bryant (University of Otago)
- ▶ BioDS Lab:

