

Distance Measures for Tumor Evolutionary Trees

Zach DiNardo^{*1}, Kiran Tomlinson^{*1}, Anna Ritz², and Layla Oesper¹

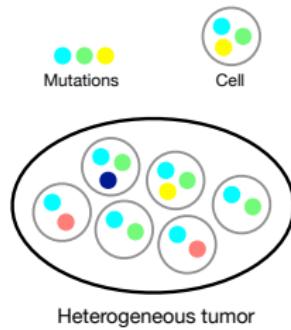
¹Department of Computer Science, Carleton College

²Department of Biology, Reed College

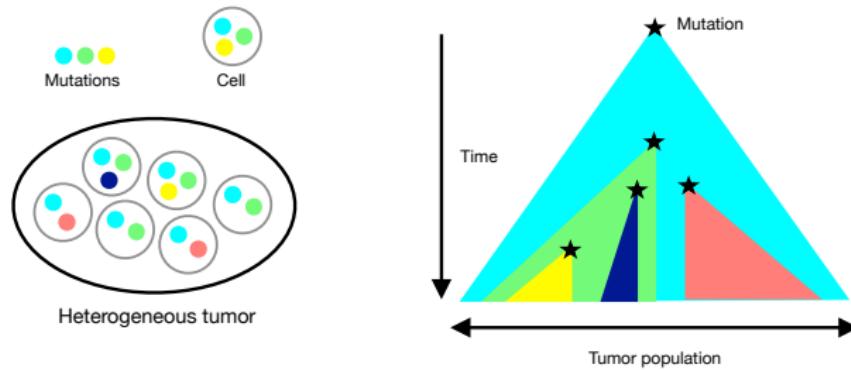
*Joint first author

May 3, 2019

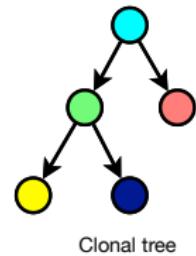
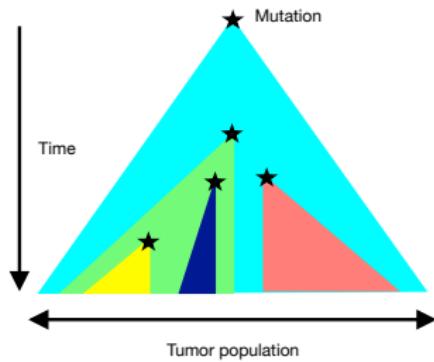
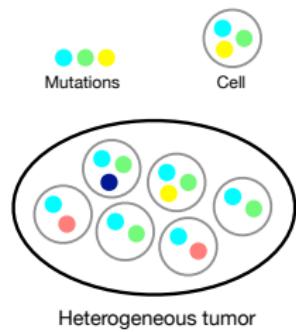
Clonal Theory (Nowell 1976)



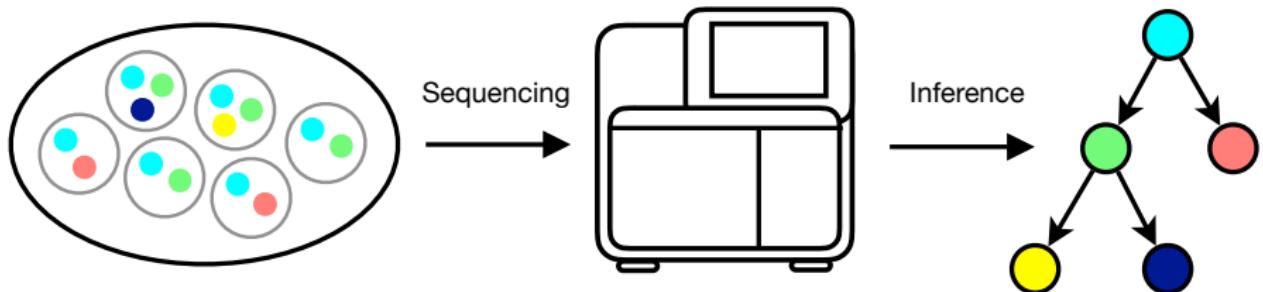
Clonal Theory (Nowell 1976)



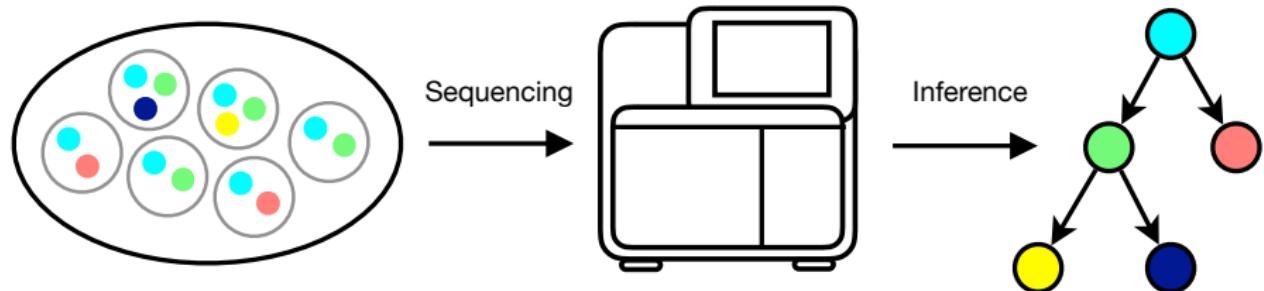
Clonal Theory (Nowell 1976)



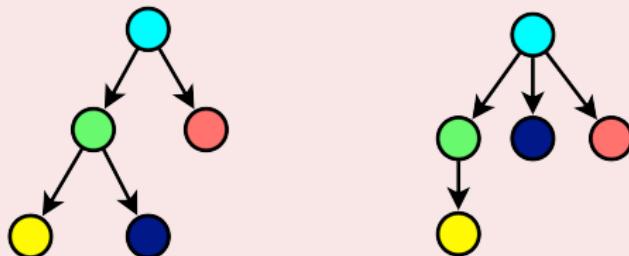
Clonal Tree Inference



Clonal Tree Inference



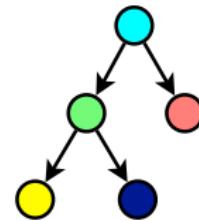
How do we compare two clonal trees?



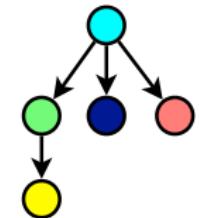
Outline

1 Introduction

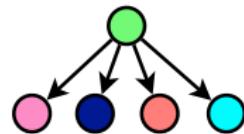
- Many Trees
- Uses for Distance Measures
- Existing Distance Measures



2 Methods

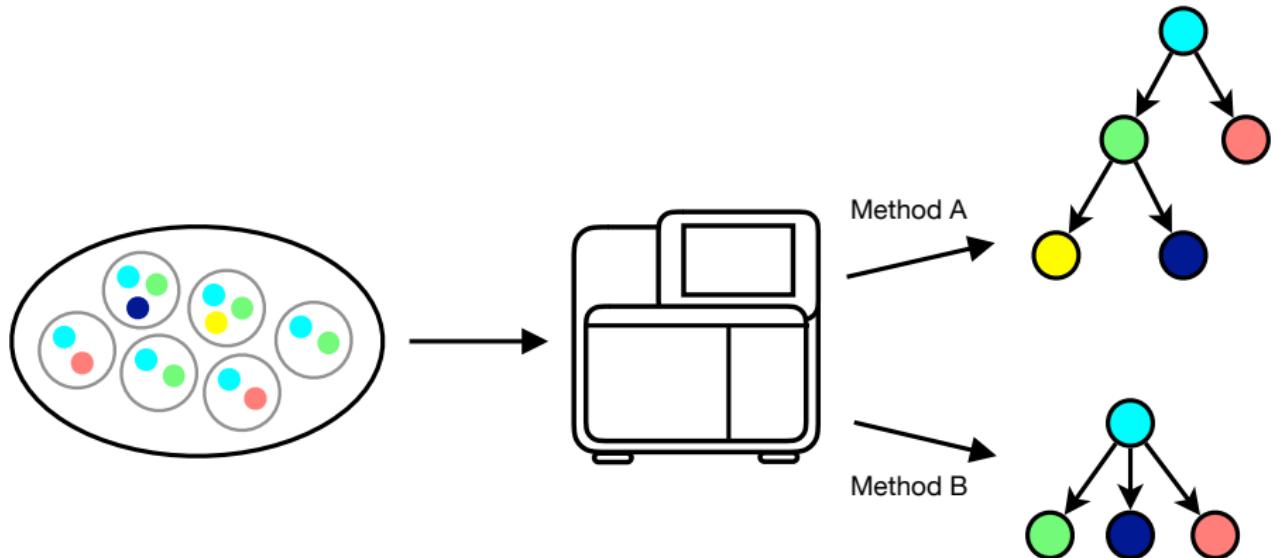


3 Results

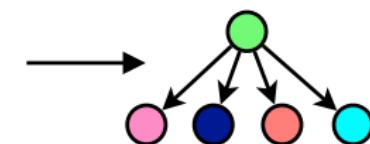
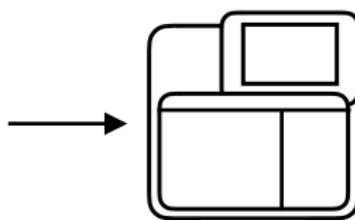
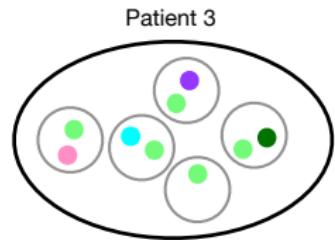
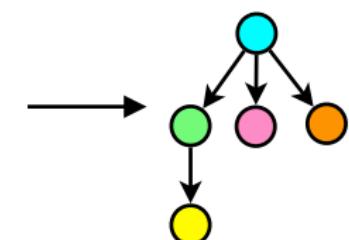
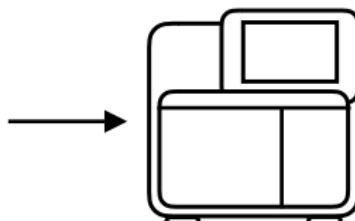
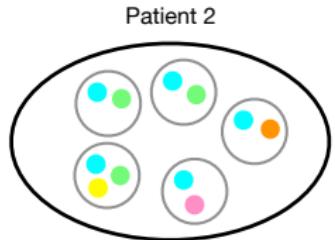
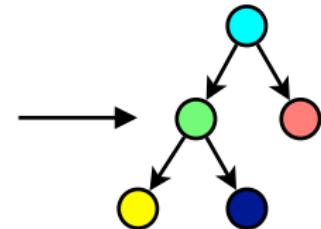
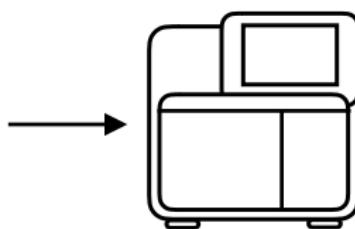
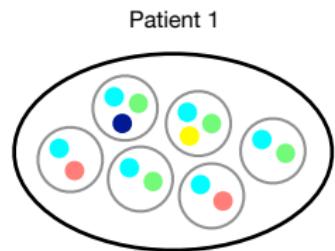


4 Conclusions

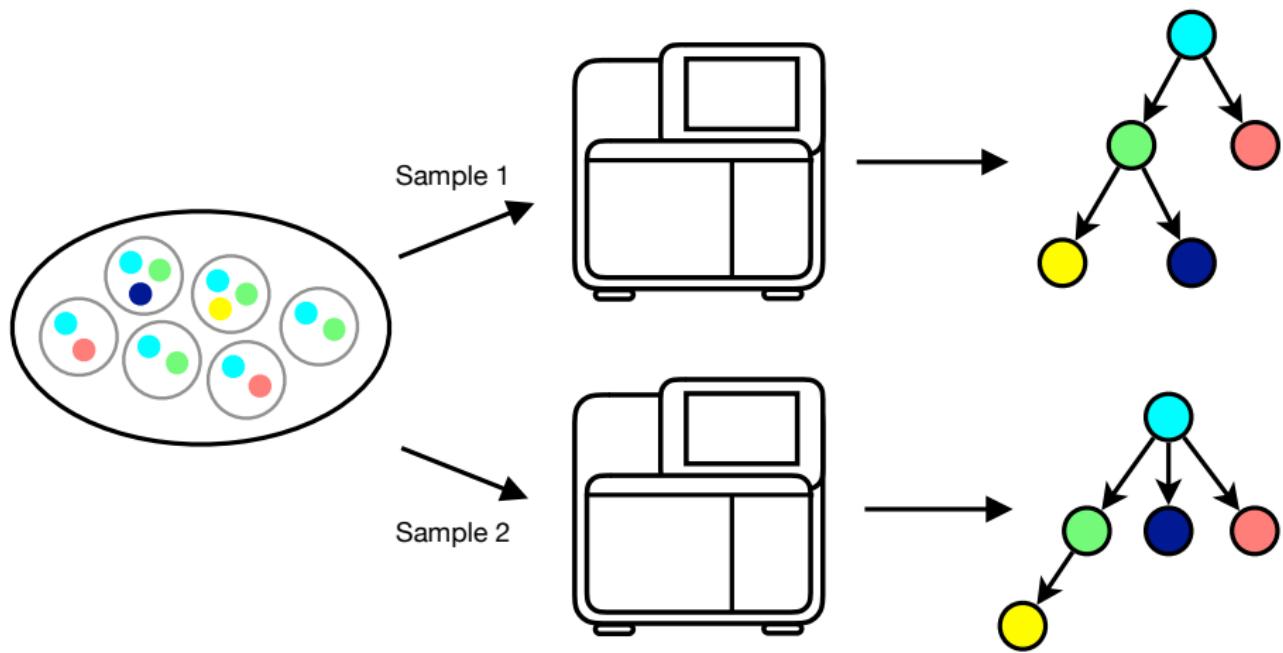
Multiple Inference Methods



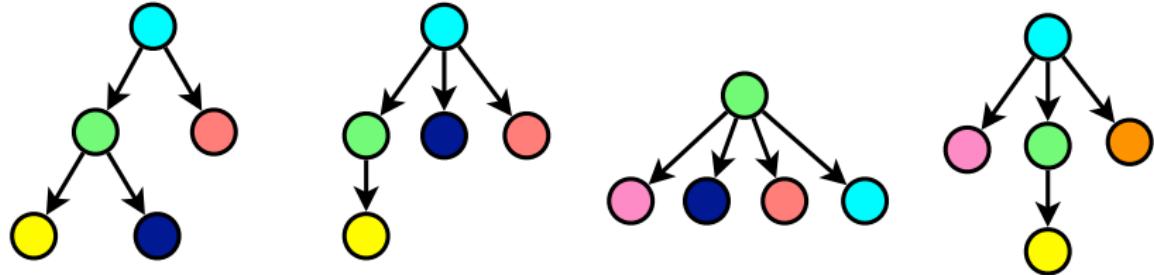
Inter-Patient Comparison



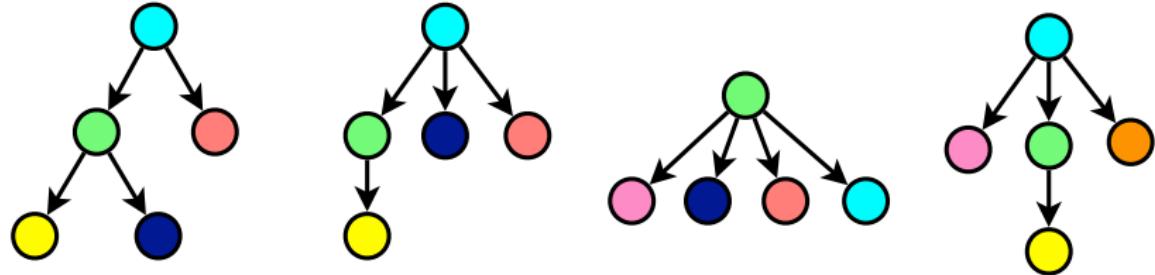
Intra-Patient Comparison



Need for Distance Measures



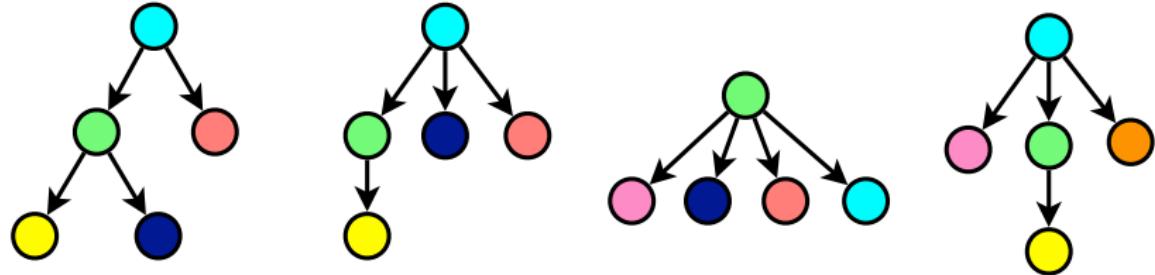
Need for Distance Measures



Uses

- ① Comparing/evaluating inferred trees

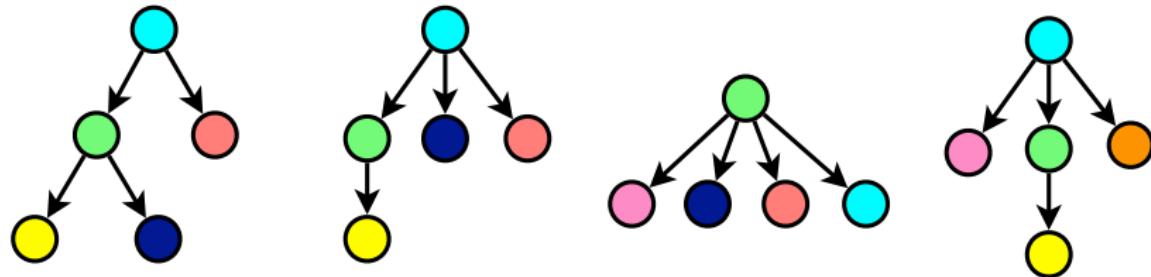
Need for Distance Measures



Uses

- ① Comparing/evaluating inferred trees
- ② Clustering trees

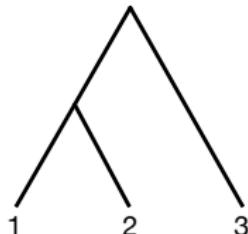
Need for Distance Measures



Uses

- ① Comparing/evaluating inferred trees
- ② Clustering trees
- ③ Inference/consensus methods (e.g., Govek et al. 2018)

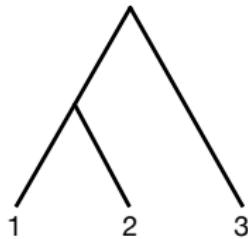
Existing Distance Measures



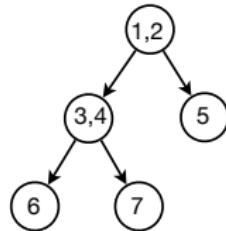
Phylogenetic trees

- ① Robinson-Foulds distance
(Robinson & Foulds 1981)
- ② Quartet distance
(Estabrook et al. 1985)
- ③ Triplet distance
(Critchlow et al. 1996)

Existing Distance Measures



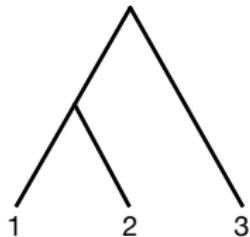
Phylogenetic trees



Clonal trees

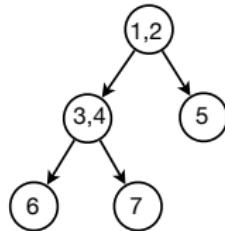
- Robinson-Foulds distance
(Robinson & Foulds 1981)
- Quartet distance
(Estabrook et al. 1985)
- Triplet distance
(Critchlow et al. 1996)

Existing Distance Measures



Phylogenetic trees

- Robinson-Foulds distance
(Robinson & Foulds 1981)
- Quartet distance
(Estabrook et al. 1985)
- Triplet distance
(Critchlow et al. 1996)



Clonal trees

- ➊ MLTED
(Karpov et al. 2018)
- ➋ A-D distance
(Govek et al. 2018)
- ➌ Rearrangement distance
(Bernardini et al. 2019)

Outline

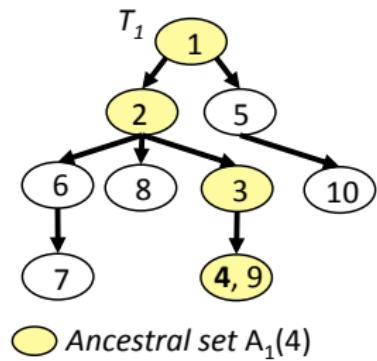
1 Introduction

2 Methods

- Definitions
- CASet
- DISC

3 Results

4 Conclusions

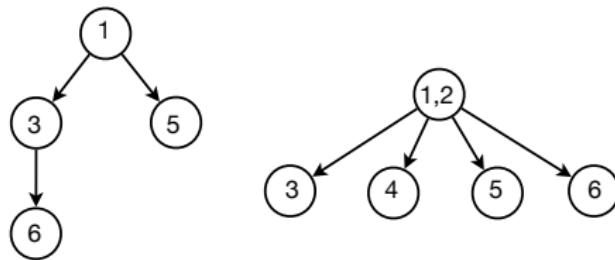


Trees

Definition

A *clonal tree* is a multi-labeled tree with unique labels.

Clonal trees:



Trees

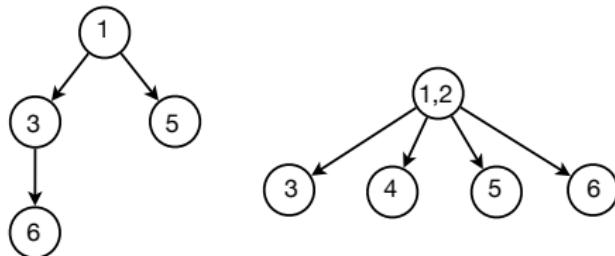
Definition

A *clonal tree* is a multi-labeled tree with unique labels.

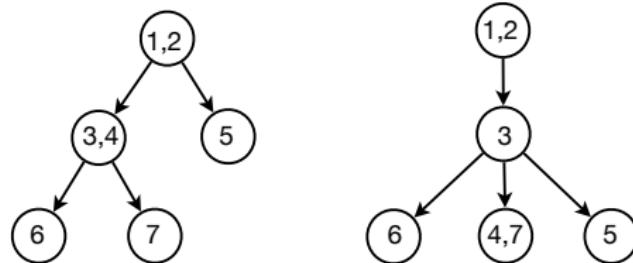
Definition

An *m-clonal tree* is a clonal tree with labels $1, \dots, m$.

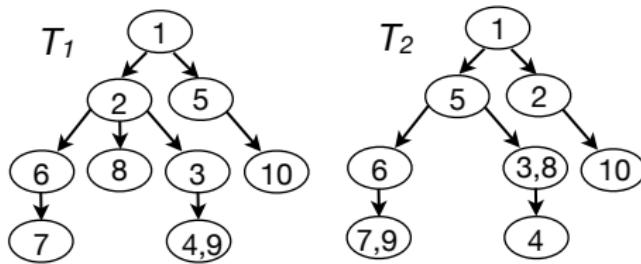
Clonal trees:



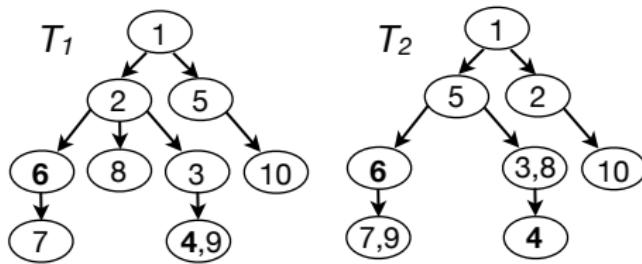
7-clonal trees:



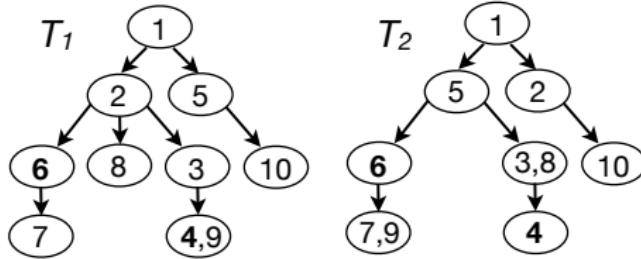
Tree Differences



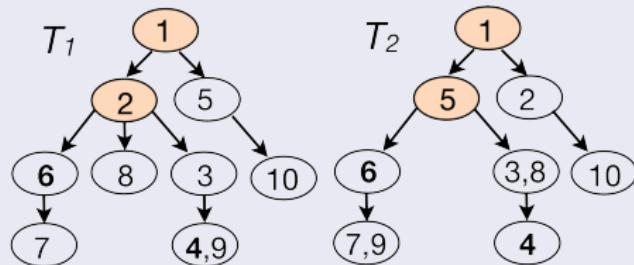
Tree Differences



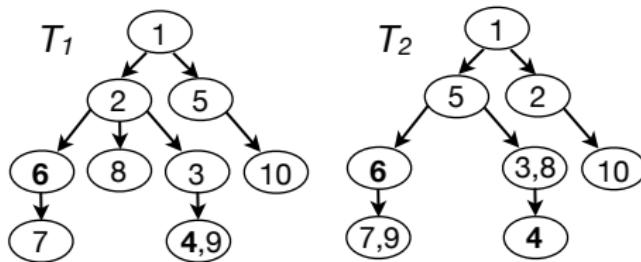
Tree Differences



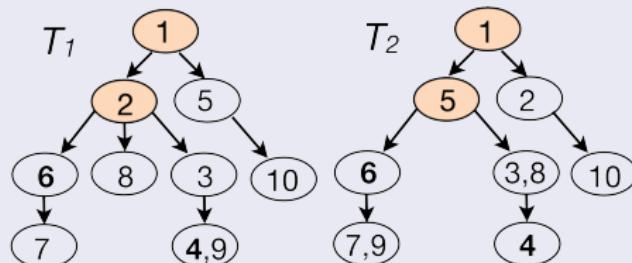
CASet



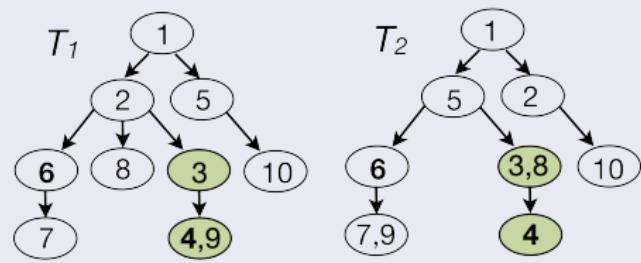
Tree Differences



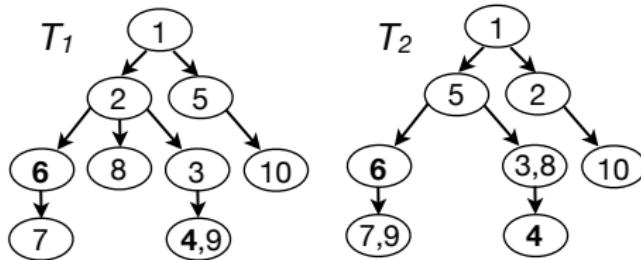
CASet



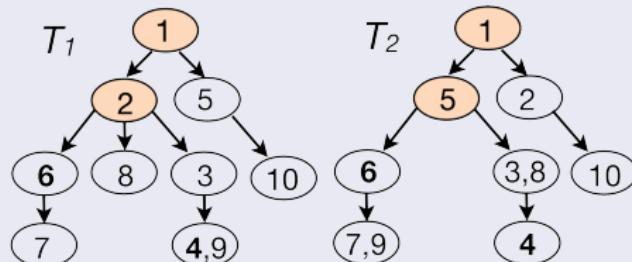
DISC



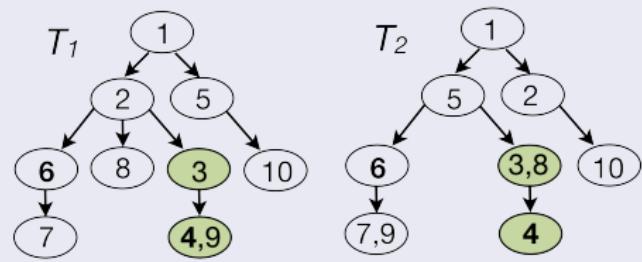
Tree Differences



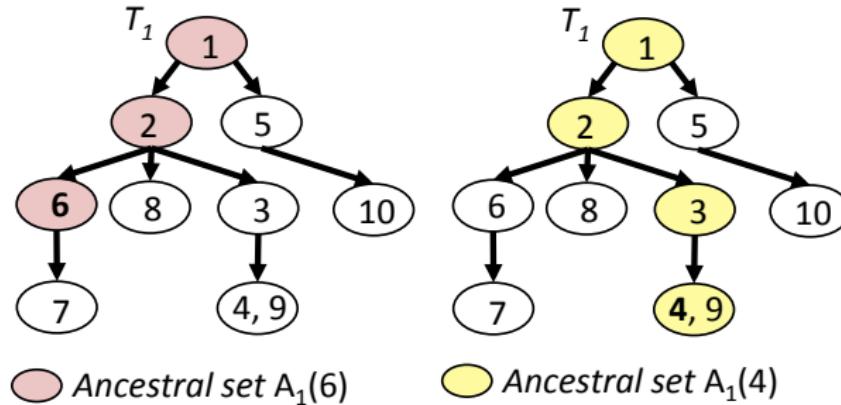
CASet



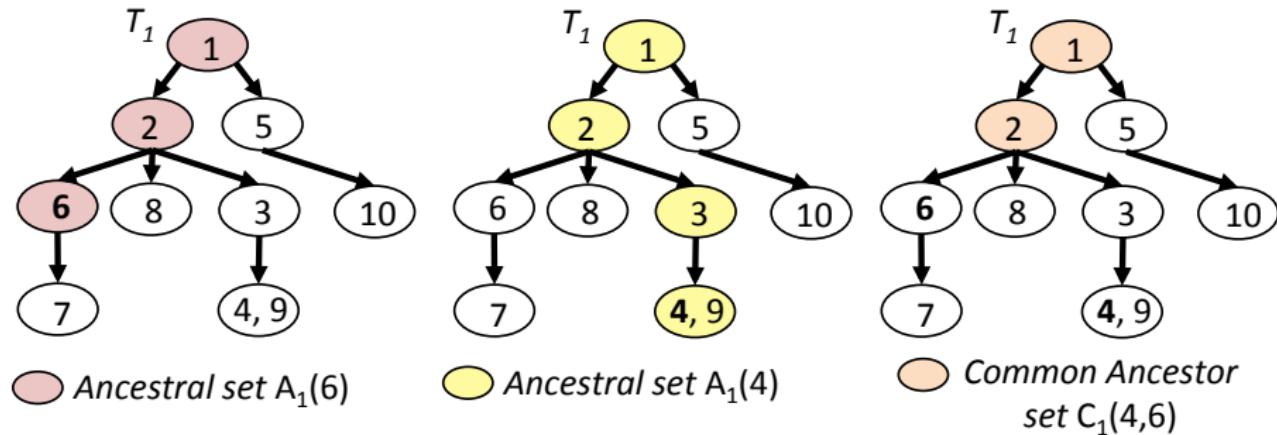
DISC



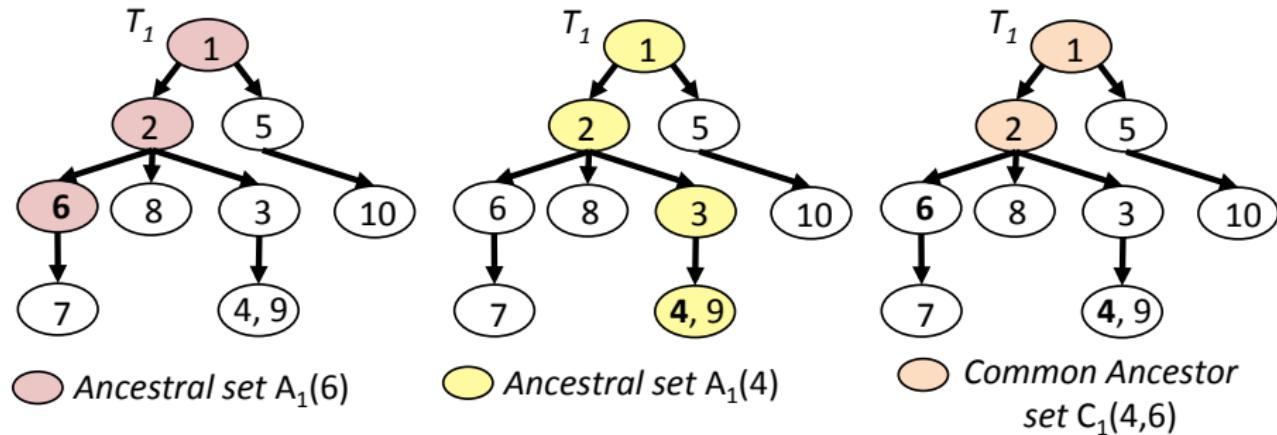
Ancestral and Common Ancestor Sets



Ancestral and Common Ancestor Sets



Ancestral and Common Ancestor Sets



Common Ancestor Set

Given a clonal tree T_k and two mutations i, j ,

$$C_k(i, j) = A_k(i) \cap A_k(j).$$

Comparing Sets

$$\text{Jacc}(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}, \quad \text{Jacc}(\emptyset, \emptyset) = 0$$

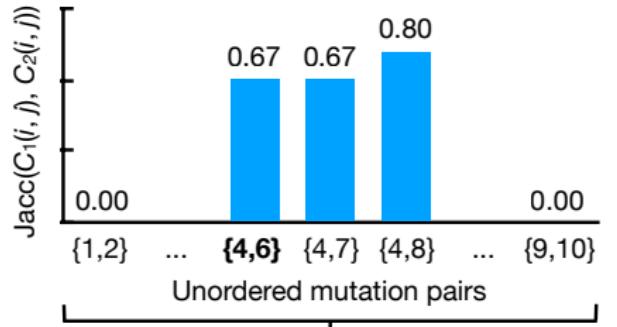
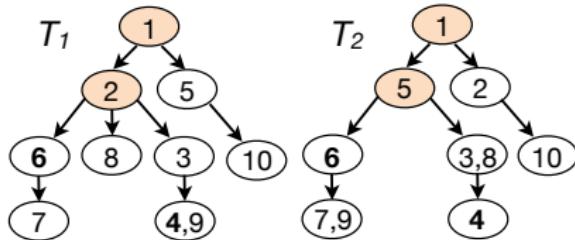
Comparing Sets

$$\text{Jacc}(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}, \quad \text{Jacc}(\emptyset, \emptyset) = 0$$

Theorem (e.g., Gilbert 1972)

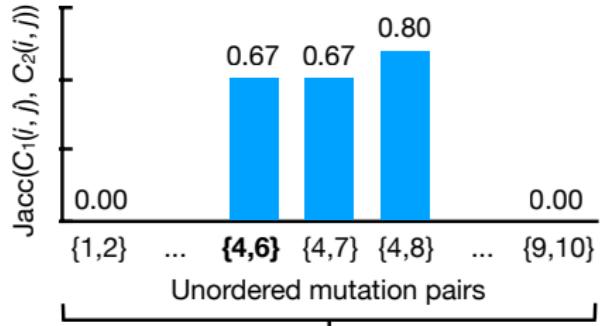
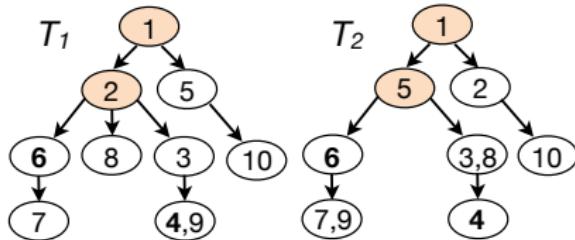
Jaccard distance is a metric on sets.

CASet



Average over all values $\rightarrow \text{CASet}(T_1, T_2) = 0.39$

CASet



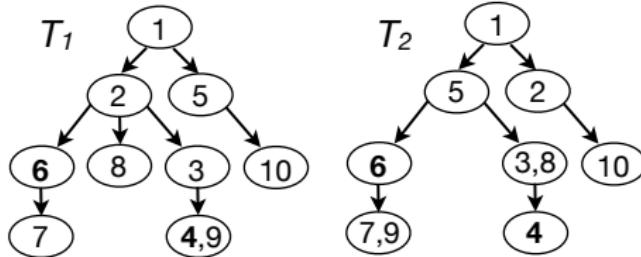
Average over all values $\rightarrow \text{CASet}(T_1, T_2) = 0.39$

CASet Distance (Common Ancestor Set)

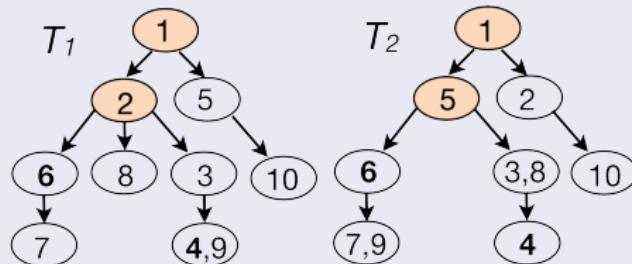
Given two m -clonal trees T_k, T_ℓ ,

$$\text{CASet}(T_k, T_\ell) = \frac{1}{\binom{m}{2}} \sum_{\{i,j\} \subseteq [m]} \text{Jacc}(C_k(i,j), C_\ell(i,j)).$$

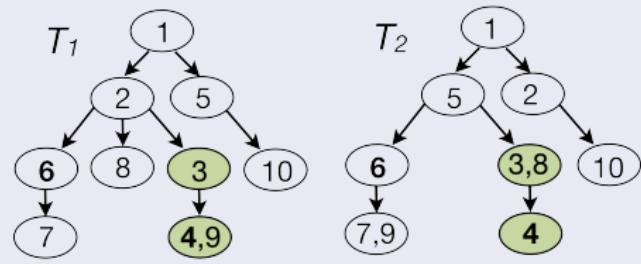
Tree Differences



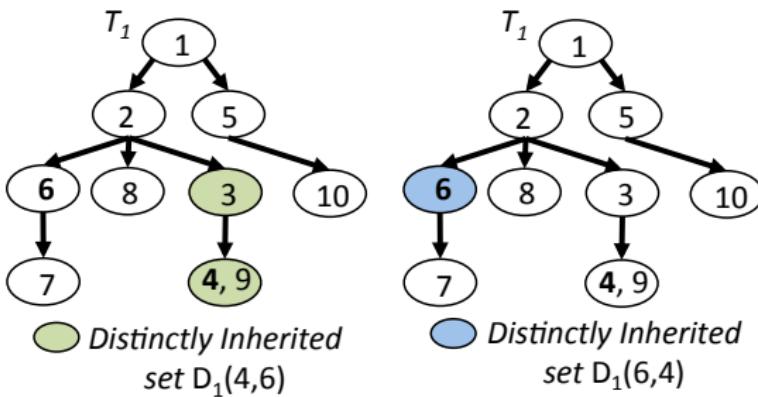
CASet



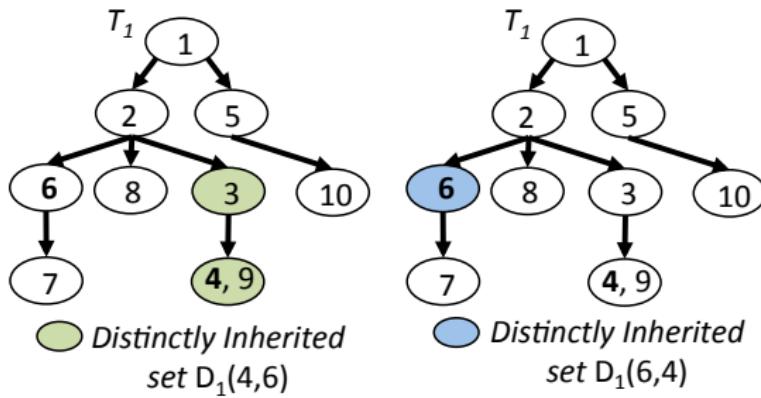
DISC



Distinctly Inherited Sets



Distinctly Inherited Sets

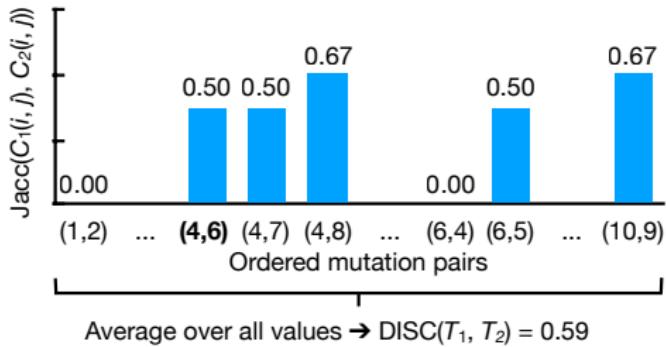
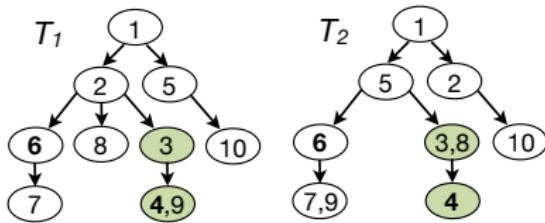


Distinctly Inherited Sets

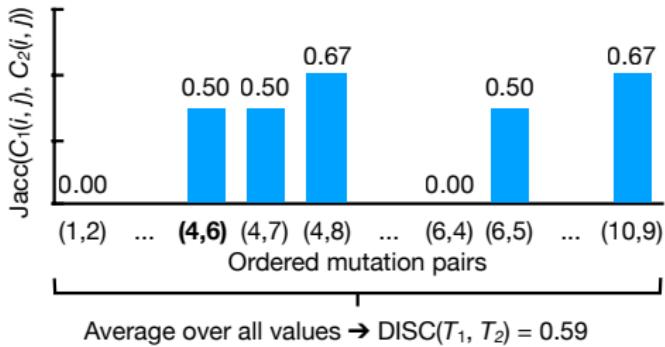
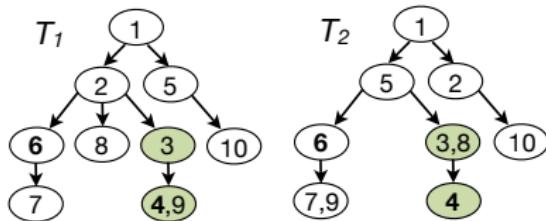
Given a clonal tree T_k and two mutations i, j ,

$$D_k(i, j) = A_k(i) \setminus A_k(j).$$

DISC



DISC



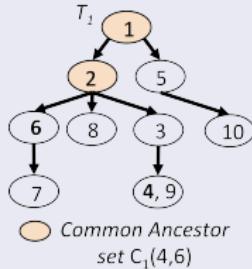
DISC Distance (Distinctly Inherited Set Comparison)

Given two m -clonal trees T_k, T_ℓ ,

$$\text{DISC}(T_k, T_\ell) = \frac{1}{m(m-1)} \sum_{\substack{(i,j) \in [m]^2 \\ i \neq j}} \text{Jacc}(D_k(i,j), D_\ell(i,j)).$$

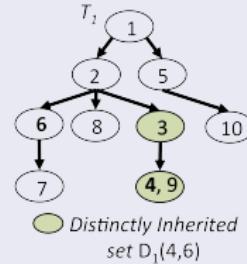
Metric Properties

CASet(T_k, T_ℓ)



$$= \frac{1}{\binom{m}{2}} \sum_{\{i,j\} \subseteq [m]} \text{Jacc}(C_k(i,j), C_\ell(i,j))$$

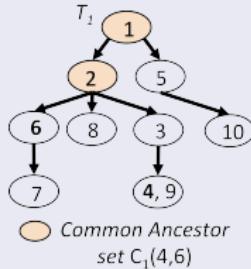
DISC(T_k, T_ℓ)



$$= \frac{1}{m(m-1)} \sum_{\substack{(i,j) \in [m]^2 \\ i \neq j}} \text{Jacc}(D_k(i,j), D_\ell(i,j))$$

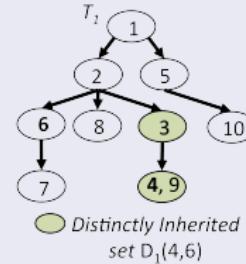
Metric Properties

CASet(T_k, T_ℓ)



$$= \frac{1}{\binom{m}{2}} \sum_{\{i,j\} \subseteq [m]} \text{Jacc}(C_k(i,j), C_\ell(i,j))$$

DISC(T_k, T_ℓ)



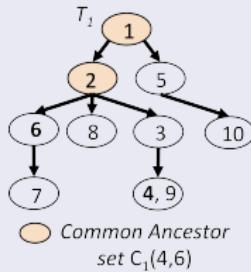
$$= \frac{1}{m(m-1)} \sum_{\substack{(i,j) \in [m]^2 \\ i \neq j}} \text{Jacc}(D_k(i,j), D_\ell(i,j))$$

Theorem

CASet and DISC are metrics on m -clonal trees.

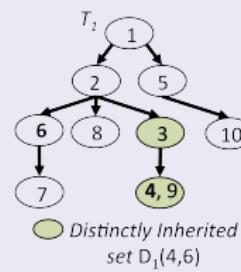
Different Label Sets: Union

CASet_U(T_k, T_ℓ)



$$= \frac{1}{\binom{|U_{k,\ell}|}{2}} \sum_{\{i,j\} \subseteq U_{k,\ell}} \text{Jacc}(C_k(i,j), C_\ell(i,j))$$

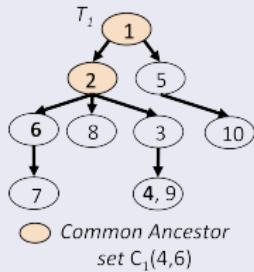
DISC_U(T_k, T_ℓ)



$$= \frac{1}{|U_{k,l}|(|U_{k,\ell}| - 1)} \sum_{\substack{(i,j) \in U_{k,\ell}^2 \\ i \neq j}} \text{Jacc}(D_k(i,j), D_\ell(i,j))$$

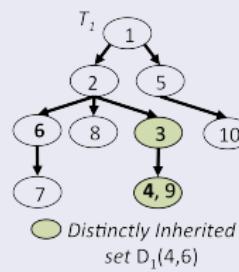
Different Label Sets: Union

CASet \cup (T_k, T_ℓ)

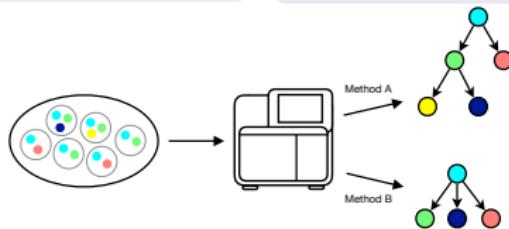


$$= \frac{1}{\binom{|U_{k,\ell}|}{2}} \sum_{\{i,j\} \subseteq U_{k,\ell}} \text{Jacc}(C_k(i,j), C_\ell(i,j))$$

DISC \cup (T_k, T_ℓ)

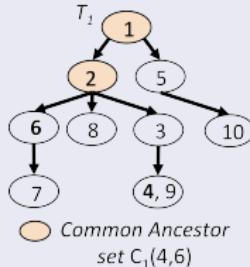


$$= \frac{1}{|U_{k,\ell}|(|U_{k,\ell}| - 1)} \sum_{\substack{(i,j) \in U_{k,\ell}^2 \\ i \neq j}} \text{Jacc}(D_k(i,j), D_\ell(i,j))$$



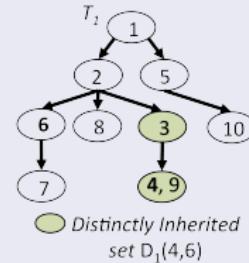
Different Label Sets: Intersection

CASet $\cap(T_k, T_\ell)$



$$= \frac{1}{\binom{|I_{k,\ell}|}{2}} \sum_{\{i,j\} \subseteq I_{k,\ell}} \text{Jacc}(C_k(i,j), C_\ell(i,j))$$

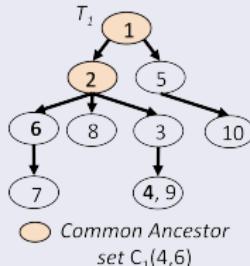
DISC $\cap(T_k, T_\ell)$



$$= \frac{1}{|I_{k,I}|(|I_{k,\ell}| - 1)} \sum_{\substack{(i,j) \in I_{k,\ell}^2 \\ i \neq j}} \text{Jacc}(D_k(i,j), D_\ell(i,j))$$

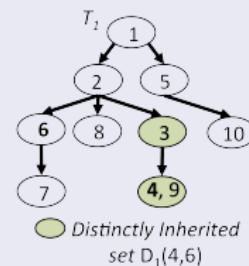
Different Label Sets: Intersection

$\text{CASet}_{\cap}(T_k, T_{\ell})$

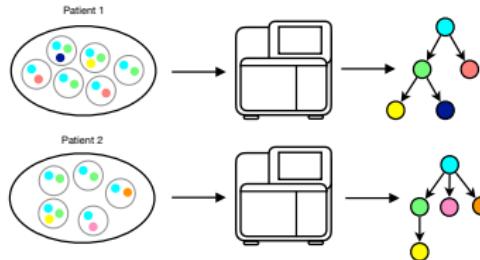


$$= \frac{1}{\binom{|I_{k,\ell}|}{2}} \sum_{\{i,j\} \subseteq I_{k,\ell}} \text{Jacc}(C_k(i,j), C_{\ell}(i,j))$$

$\text{DISC}_{\cap}(T_k, T_{\ell})$

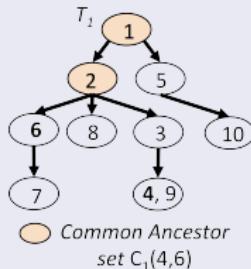


$$= \frac{1}{|I_{k,\ell}|(|I_{k,\ell}| - 1)} \sum_{\substack{(i,j) \in I_{k,\ell}^2 \\ i \neq j}} \text{Jacc}(D_k(i,j), D_{\ell}(i,j))$$



Overview of Distance Measures

CASet(T_k, T_ℓ)



$$= \frac{1}{\binom{m}{2}} \sum_{\{i,j\} \subseteq [m]} \text{Jacc}(C_k(i,j), C_\ell(i,j))$$

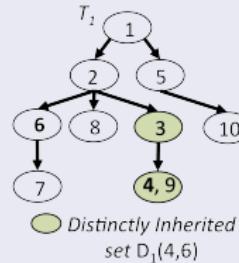
Union

$$\sum_{\{i,j\} \subseteq U_{k,\ell}}$$

Intersection

$$\sum_{\{i,j\} \subseteq I_{k,\ell}}$$

DISC(T_k, T_ℓ)



$$= \frac{1}{m(m-1)} \sum_{\substack{(i,j) \in [m]^2 \\ i \neq j}} \text{Jacc}(D_k(i,j), D_\ell(i,j))$$

Union

$$\sum_{\substack{(i,j) \in U_{k,\ell}^2 \\ i \neq j}}$$

Intersection

$$\sum_{\substack{(i,j) \in I_{k,\ell}^2 \\ i \neq j}}$$

Outline

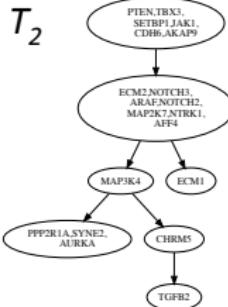
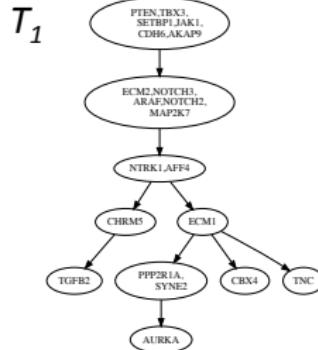
1 Introduction

2 Methods

3 Results

- Simulated Data
- Real Data

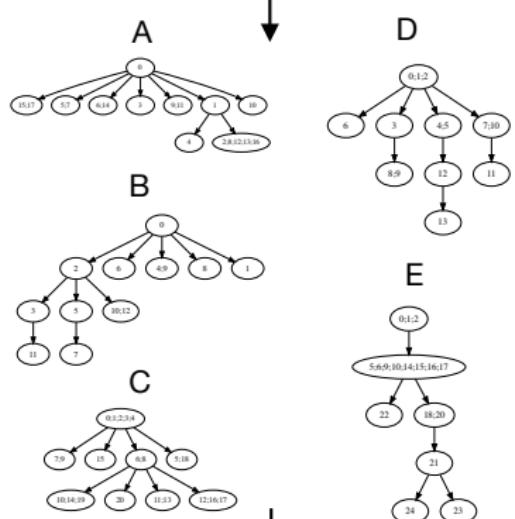
4 Conclusions



Simulated Data

OncoLib
(El-Kebir Group 2018)

Simulate tumor evolution

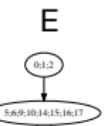
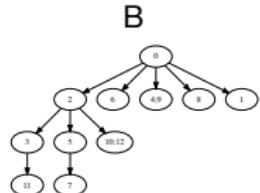
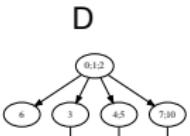
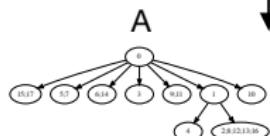


Simulate sequencing

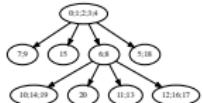
Simulated Data

OncoLib
(El-Kebir Group 2018)

Simulate tumor evolution



C



Simulate sequencing

Modified AncesTree
(El-Kebir et al. 2015, Tomlinson & Oesper 2018)

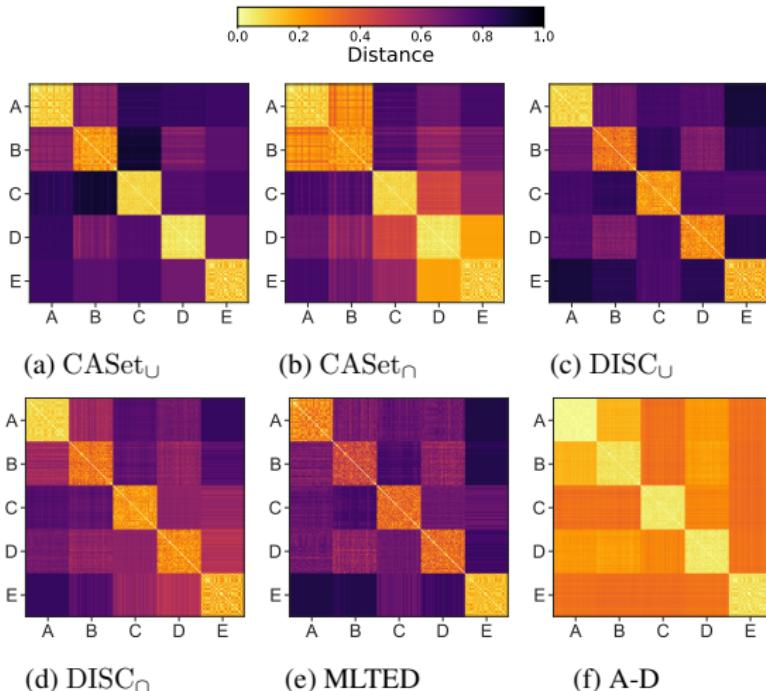
Simulated sequencing data

Enumerate compatible trees

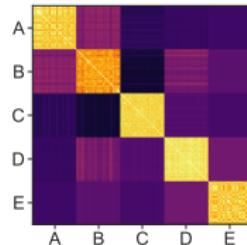
Random sampling

- Family A { ... }
- B { ... }
- C { ... } (50 trees each)
- D { ... }
- E { ... }

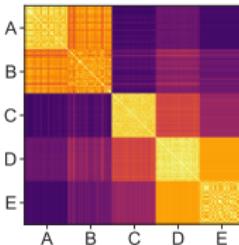
Clustering OncoLib Trees



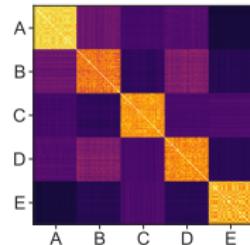
Clustering OncoLib Trees



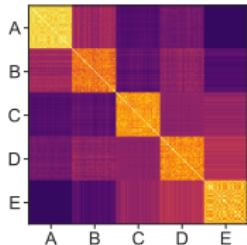
(a) CASet_U (0.81)



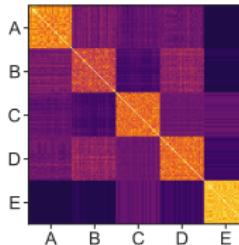
(b) CASet_∩ (0.57)



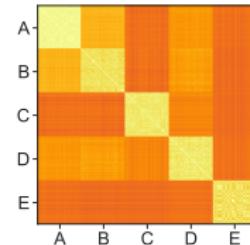
(c) DISC_U (0.70)



(d) DISC_∩ (0.62)

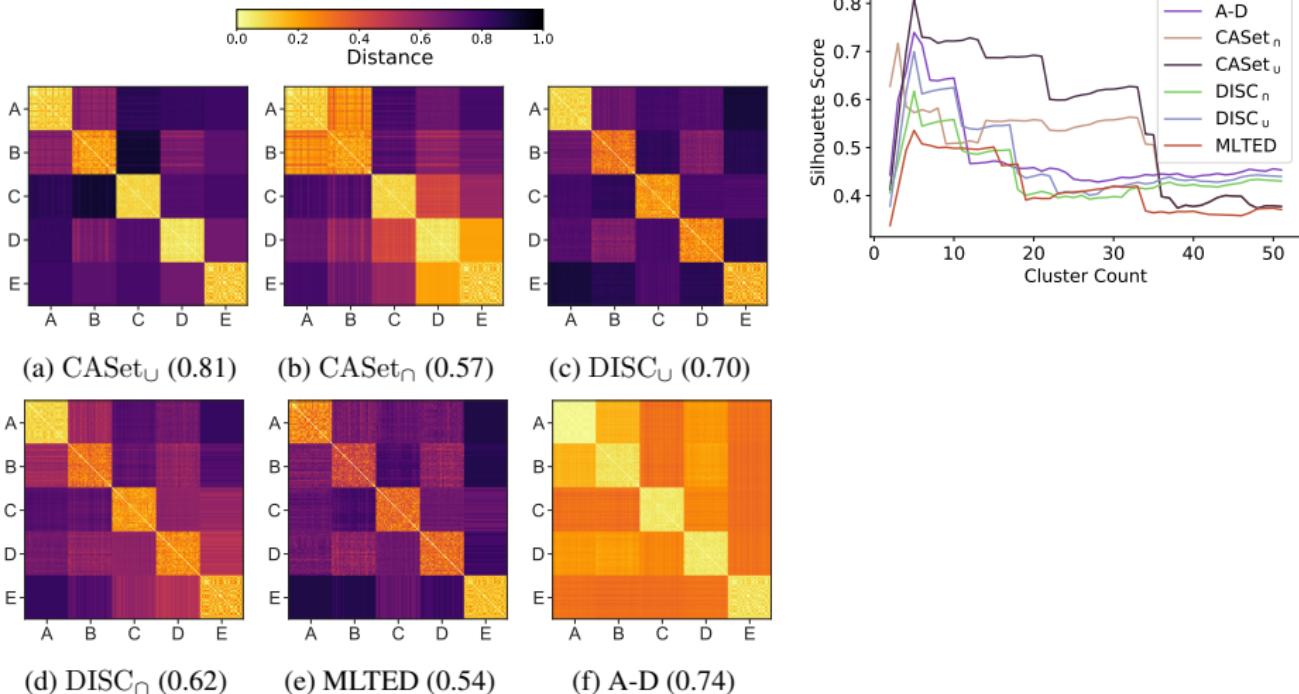


(e) MLTED (0.54)

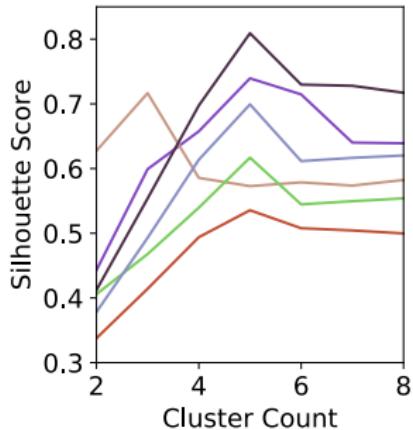
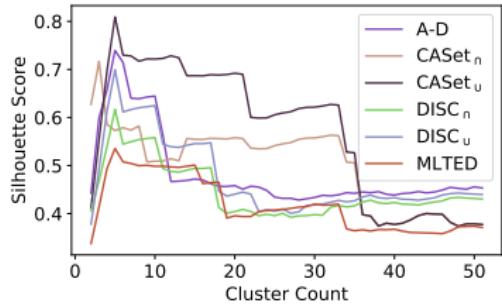
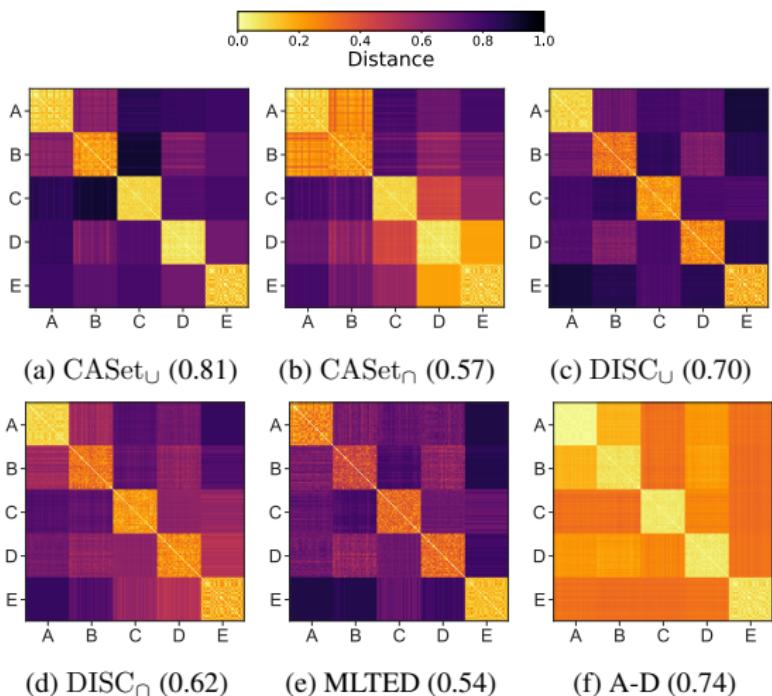


(f) A-D (0.74)

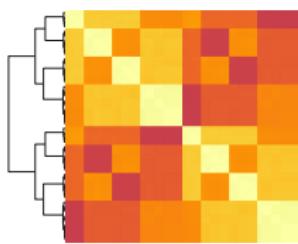
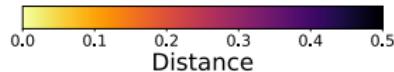
Clustering OncoLib Trees



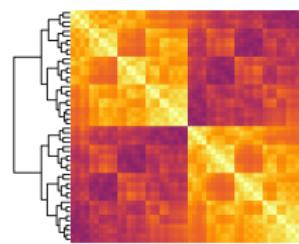
Clustering OncoLib Trees



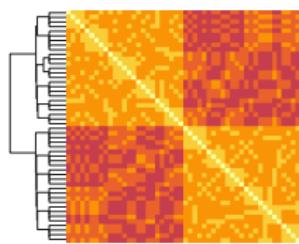
Intra-Family Structure



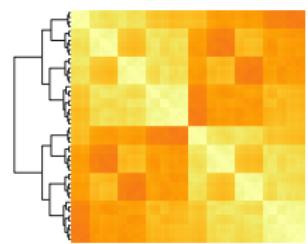
(a) CASet



(b) DISC



(c) MLTED



(d) A-D

Family E

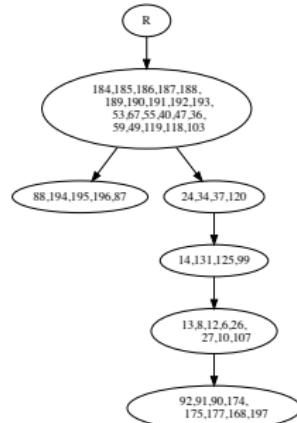
Real Datasets

① Triple negative breast cancer (Wang et al. 2014)

- Single-cell seq. at $72\times$ coverage
- Bulk deep seq. at $118,743\times$ coverage

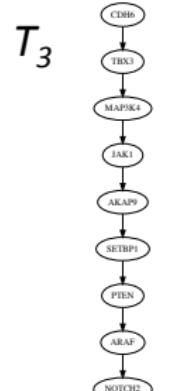
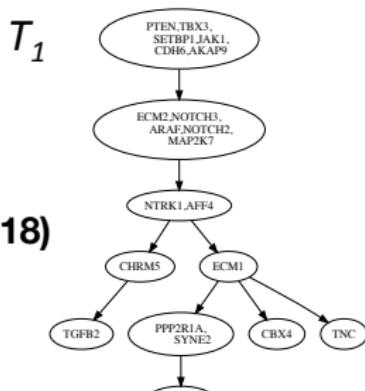
② Breast cancer xenograft (Eirew et al. 2015)

- Whole-genome seq. at $35-72\times$ coverage
- MiSeq targeted deep seq.



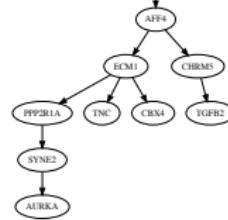
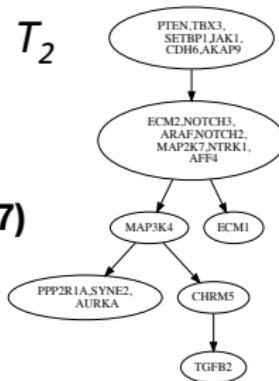
TNBC Tree Inference Analysis

PHiSCS
(Malikic et al., 2018)

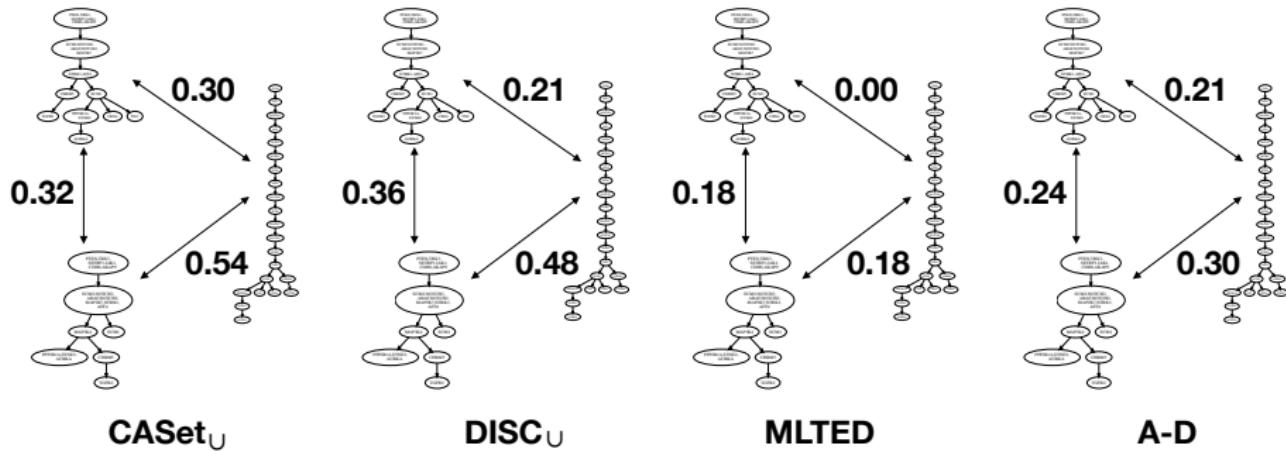


SCITE
(Jahn et al., 2016)

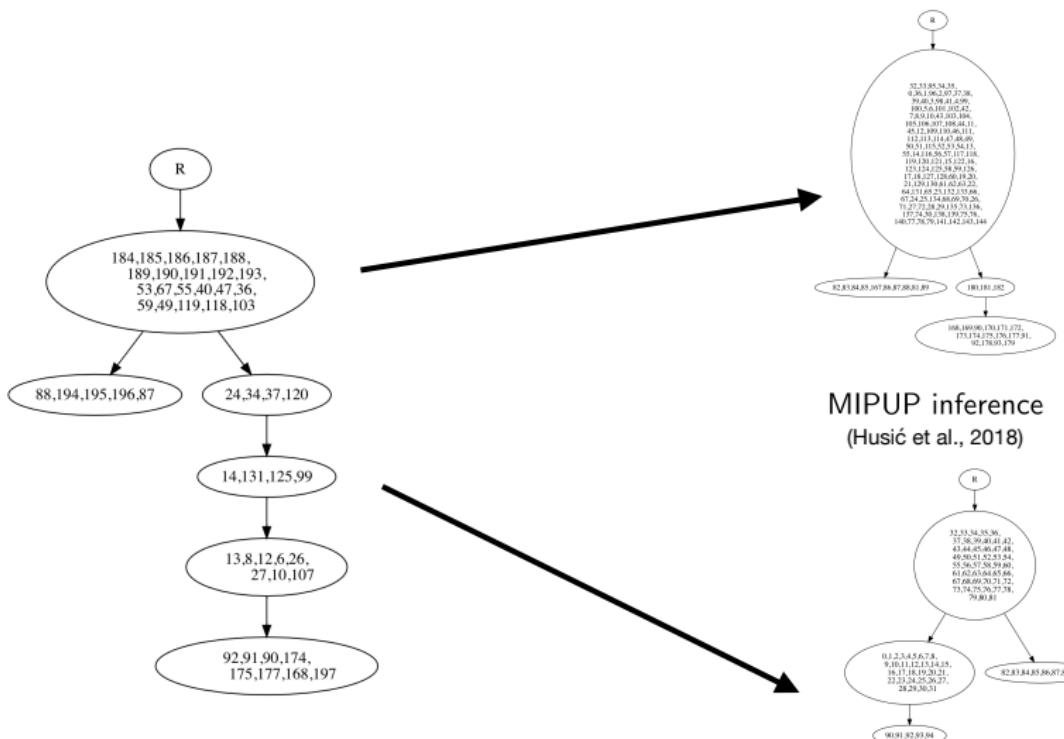
SiFit
(Zafar et al., 2017)



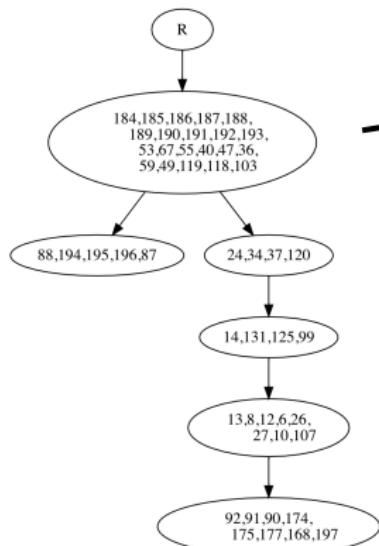
TNBC Tree Inference Analysis



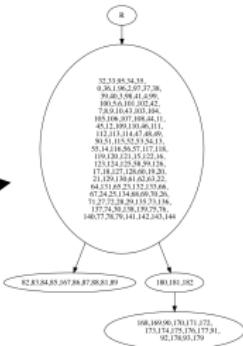
Xenoengraftment Tree Inference Analysis



Xenoengraftment Tree Inference Analysis

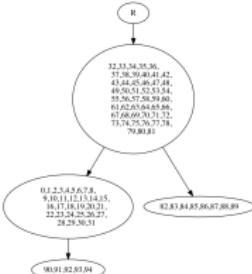


Measure	Distance
CASet _U	0.88
CASet _n	0.84
DISC _U	0.40
DISC _n	0.40
MLTED	0.80
A-D	0.70



MIPUP inference

Measure	Distance
CASet _U	0.74
CASet _n	0.78
DISC _U	0.60
DISC _n	0.38
MLTED	0.81
A-D	0.46

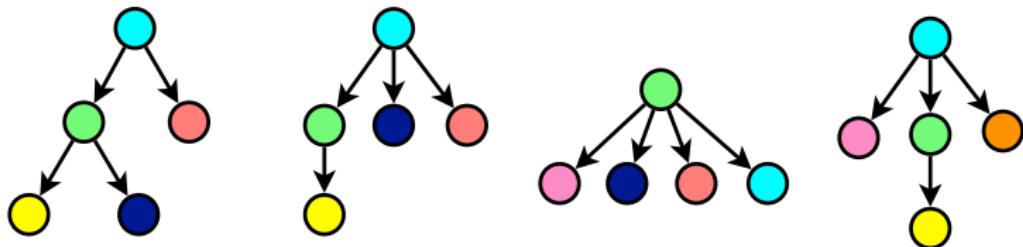


LICHeE inference

SA501 Base Tree

Conclusions

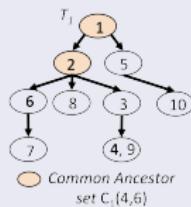
- ① Distance measures are important for tumor tree analysis
- ② We introduced two novel distance metrics, CASet and DISC
- ③ CASet_\cup clusters trees more clearly than existing measures
- ④ CASet and DISC have high resolution on simulated and real data



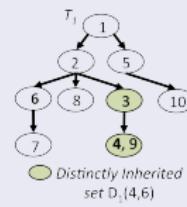
Conclusions

- ① Distance measures are important for tumor tree analysis
- ② We introduced two novel distance metrics, CASet and DISC
- ③ CASet_\cup clusters trees more clearly than existing measures
- ④ CASet and DISC have high resolution on simulated and real data

$\text{CASet}(T_k, T_\ell)$



$\text{DISC}(T_k, T_\ell)$



Union

$$\sum_{\{i,j\} \subseteq U_{k,\ell}}$$

Intersection

$$\sum_{\{i,j\} \subseteq I_{k,\ell}}$$

Union

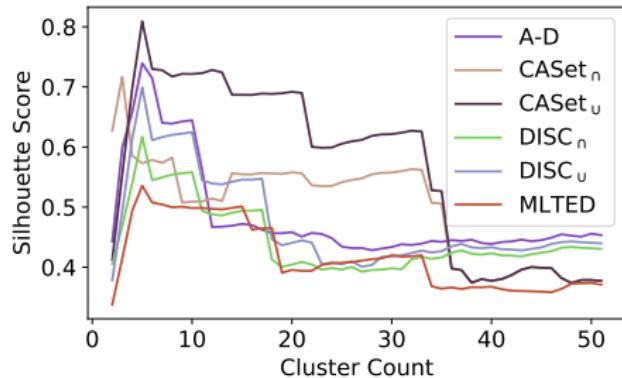
$$\sum_{\substack{(i,j) \in U_{k,\ell}^2 \\ i \neq j}}$$

Intersection

$$\sum_{\substack{(i,j) \in I_{k,\ell}^2 \\ i \neq j}}$$

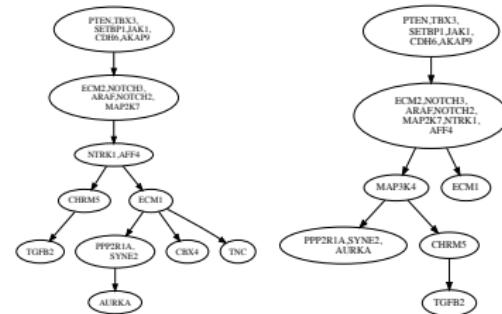
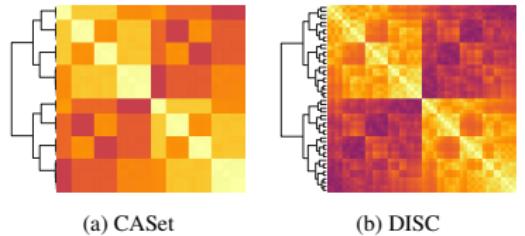
Conclusions

- ① Distance measures are important for tumor tree analysis
- ② We introduced two novel distance metrics, CASet and DISC
- ③ CASet_U clusters trees more clearly than existing measures
- ④ CASet and DISC have high resolution on simulated and real data



Conclusions

- ① Distance measures are important for tumor tree analysis
- ② We introduced two novel distance metrics, CASet and DISC
- ③ CASet_\cup clusters trees more clearly than existing measures
- ④ CASet and DISC have high resolution on simulated and real data



Acknowledgment



- This project is supported by the NSF and Elledge, Eugster, and Class of '49 Fellowships from Carleton College (to LO).
- Thank you to the Carleton College Computer Science Department, Layla Oesper, Anna Ritz, Rosa Zhou, and Thais Del Rosario Hernandez.

Availability

Implementations of CASet and DISC are available at

<https://bitbucket.org/oesperlab/stereodist>.

Preprint available at

<https://www.biorxiv.org/>.