

Implications of Non-uniqueness of Solutions in Cancer Phylogenetics

Mohammed El-Kebir

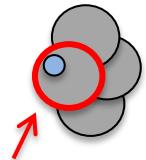
University of Illinois at Urbana Champaign,
Department of Computer Science

DSW 2019



Tumorigenesis: Cell Mutation

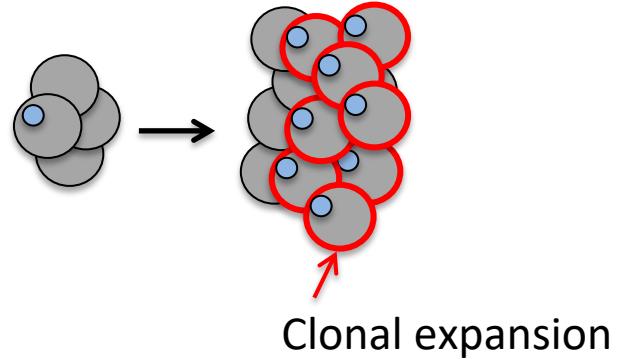
Clonal Evolution Theory of Cancer
[Nowell, 1976]



Founder
tumor cell
with somatic mutation: ●
(e.g. BRAF V600E)

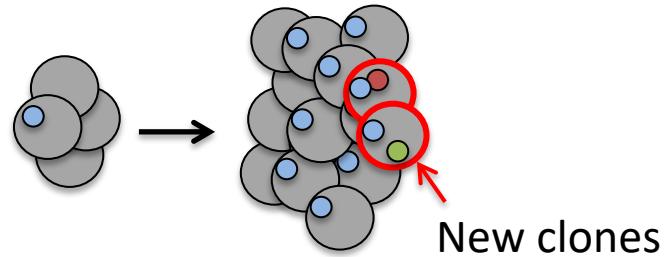
Tumorigenesis: Cell Mutation

Clonal Evolution Theory of Cancer
[Nowell, 1976]



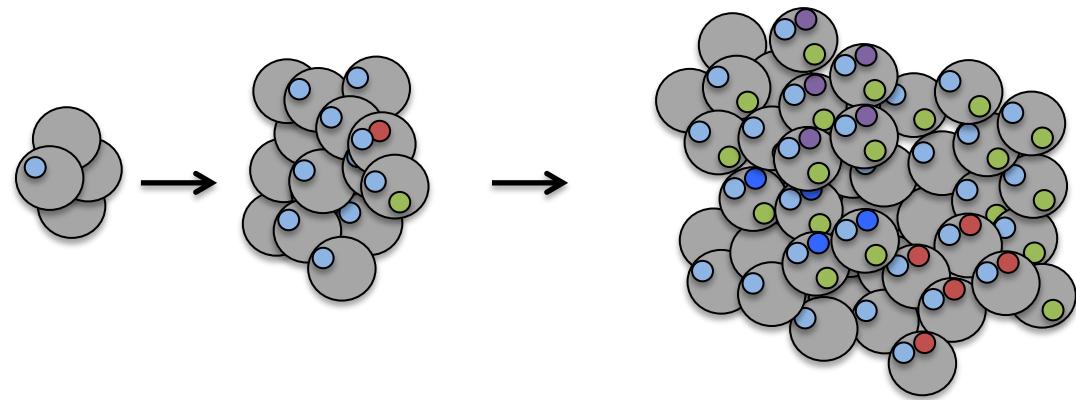
Tumorigenesis: Cell Mutation

Clonal Evolution Theory of Cancer
[Nowell, 1976]



Tumorigenesis: Cell Mutation & Division

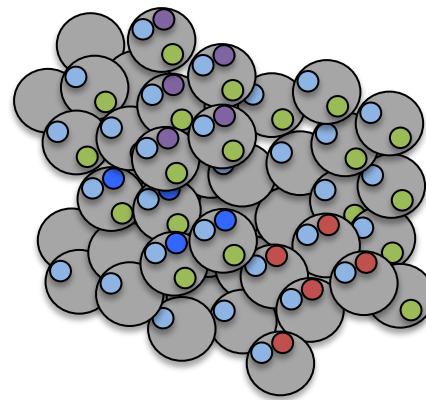
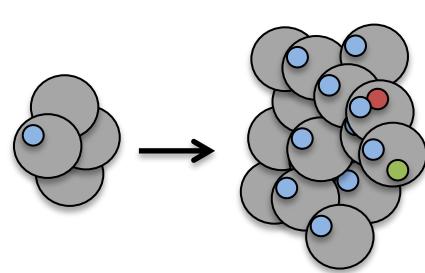
Clonal Evolution Theory of Cancer
[Nowell, 1976]



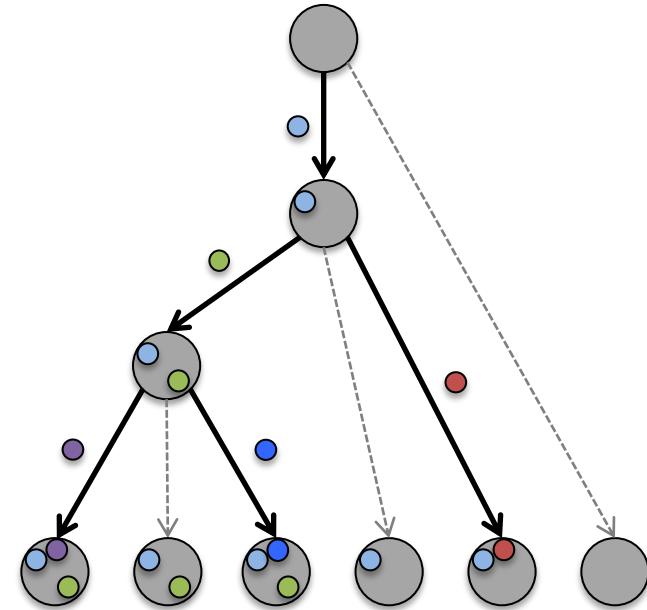
Intra-Tumor
Heterogeneity

Tumorigenesis: Cell Mutation & Division

Clonal Evolution Theory of Cancer
[Nowell, 1976]



Intra-Tumor
Heterogeneity

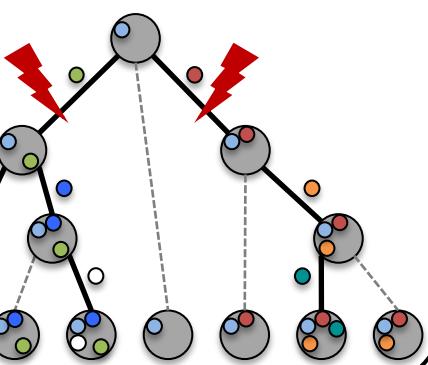


Phylogenetic
Tree T

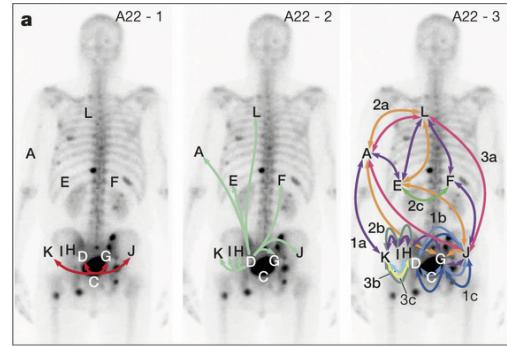
Question: Why are tumor phylogenies important?

Phylogenies are Key to Understanding Cancer

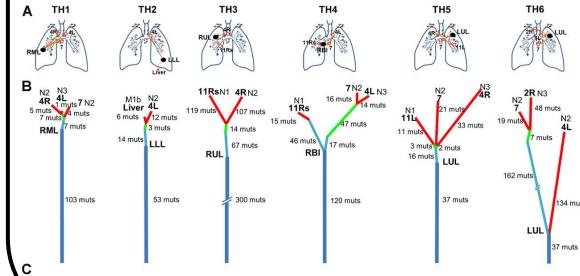
Identify targets for treatment



Understand metastatic development

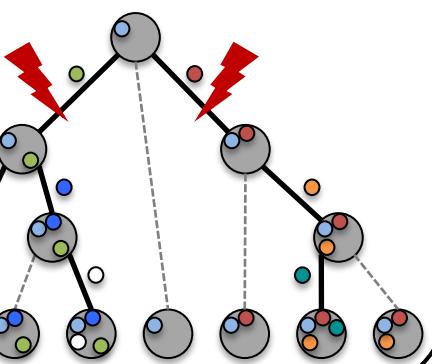


Recognize common patterns of tumor evolution across patients

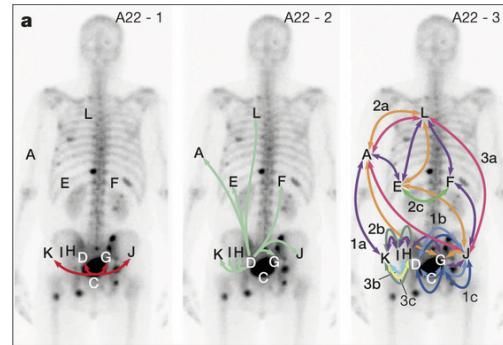


Phylogenies are Key to Understanding Cancer

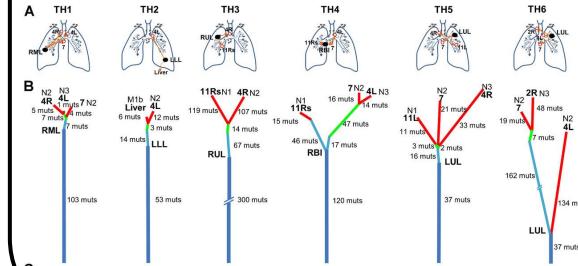
Identify targets for treatment



Understand metastatic development



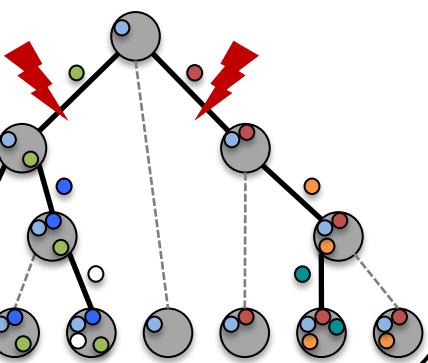
Recognize common patterns of tumor evolution across patients



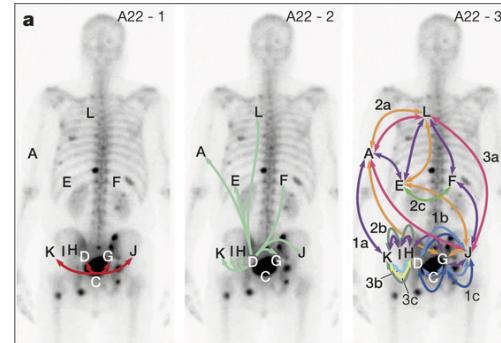
These downstream analyses **critically rely** on accurate tumor phylogeny inference

Phylogenies are Key to Understanding Cancer

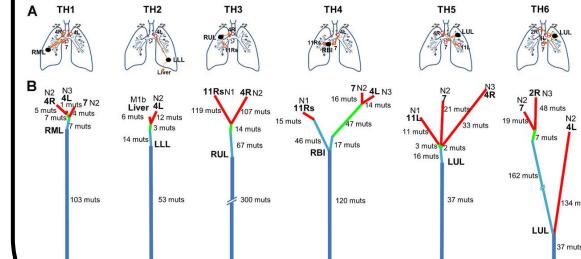
Identify targets for treatment



Understand metastatic development



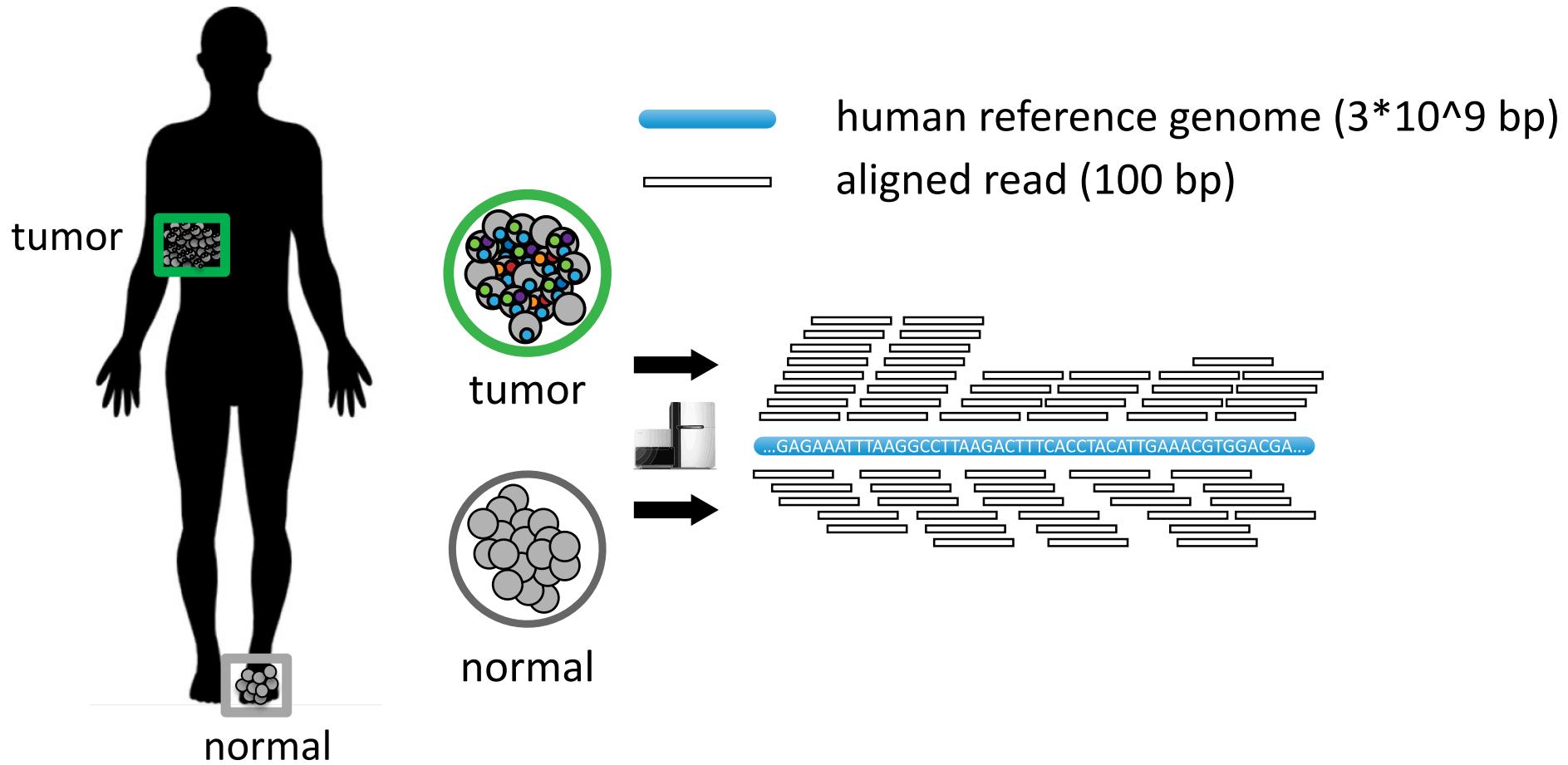
Recognize common patterns of tumor evolution across patients



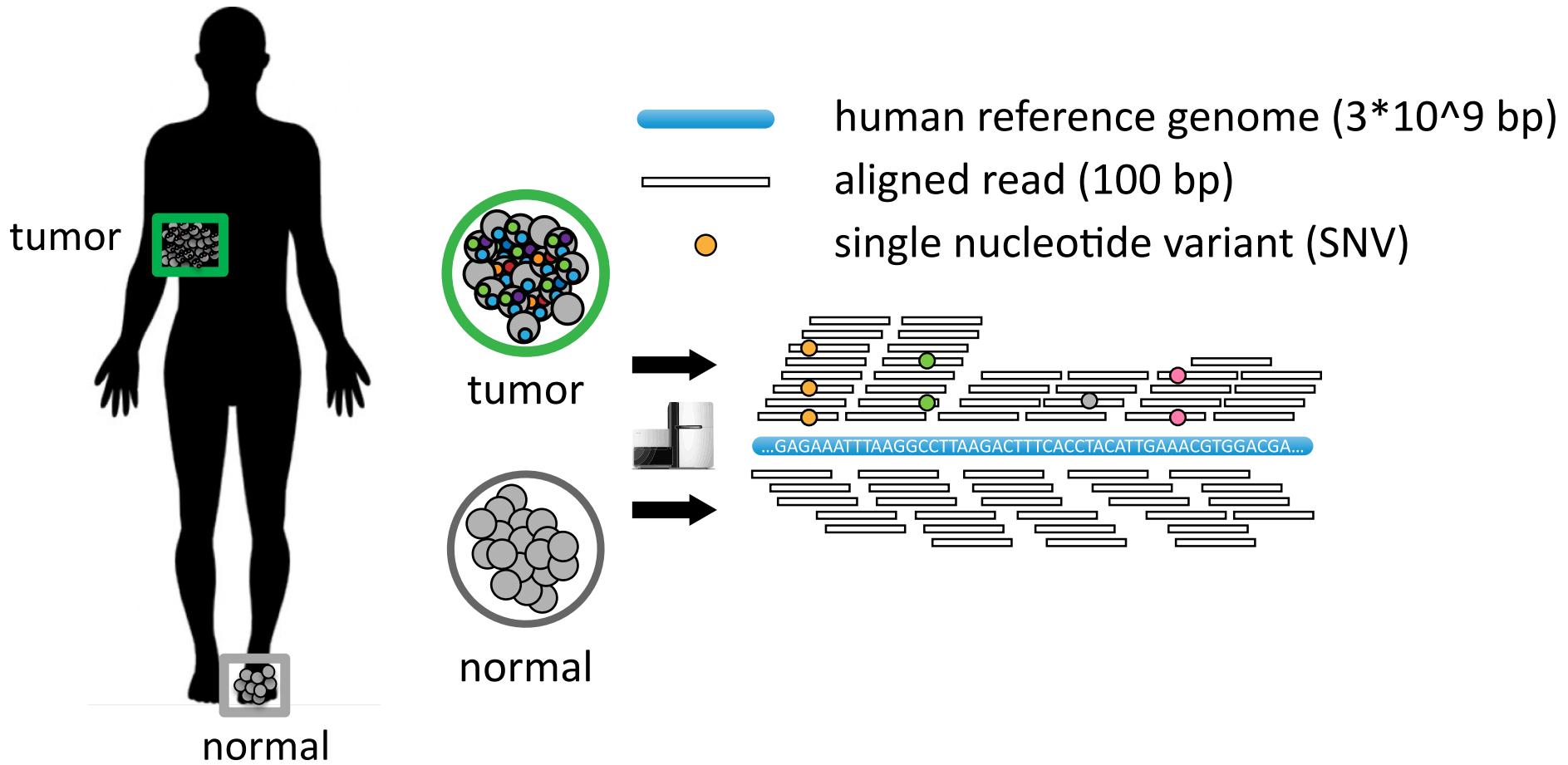
These downstream analyses **critically rely** on accurate tumor phylogeny inference

Key challenge in phylogenetics:
Accurate phylogeny inference from data at present time

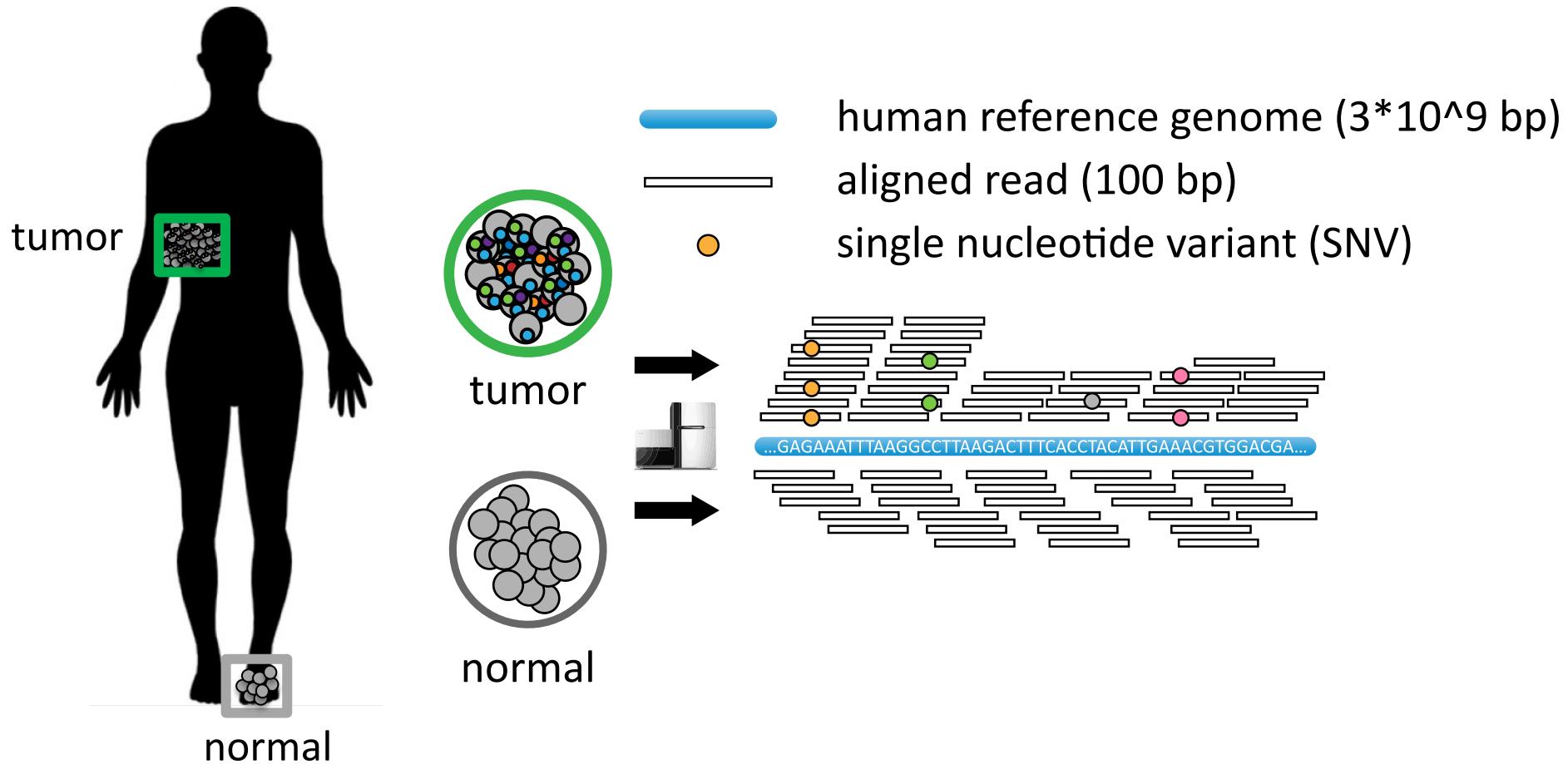
Additional Challenge in Cancer Phylogenetics



Additional Challenge in Cancer Phylogenetics



Additional Challenge in Cancer Phylogenetics



Additional challenge in cancer phylogenetics:
Phylogeny inference from **mixed bulk samples** at present time

Outline

1. Background and theory: [RECOMB-CG 2018]

- Perfect Phylogeny Mixture (PPM) problem
- Combinatorial characterization of solutions
- #PPM: exact counting and uniform sampling

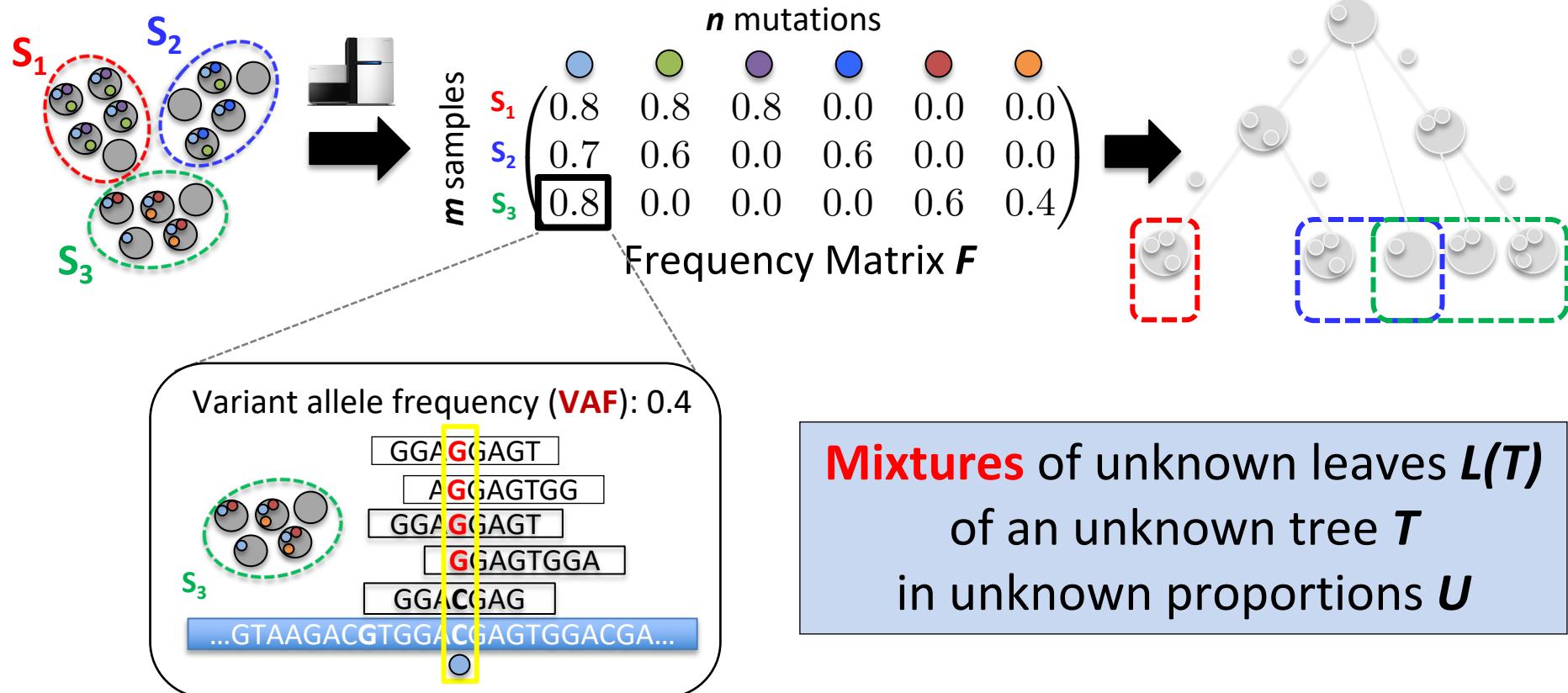
2. Simulation results: [RECOMB-CG 2018]

- What contributes to non-uniqueness?
- How to reduce non-uniqueness?
- How does non-uniqueness affect current methods?

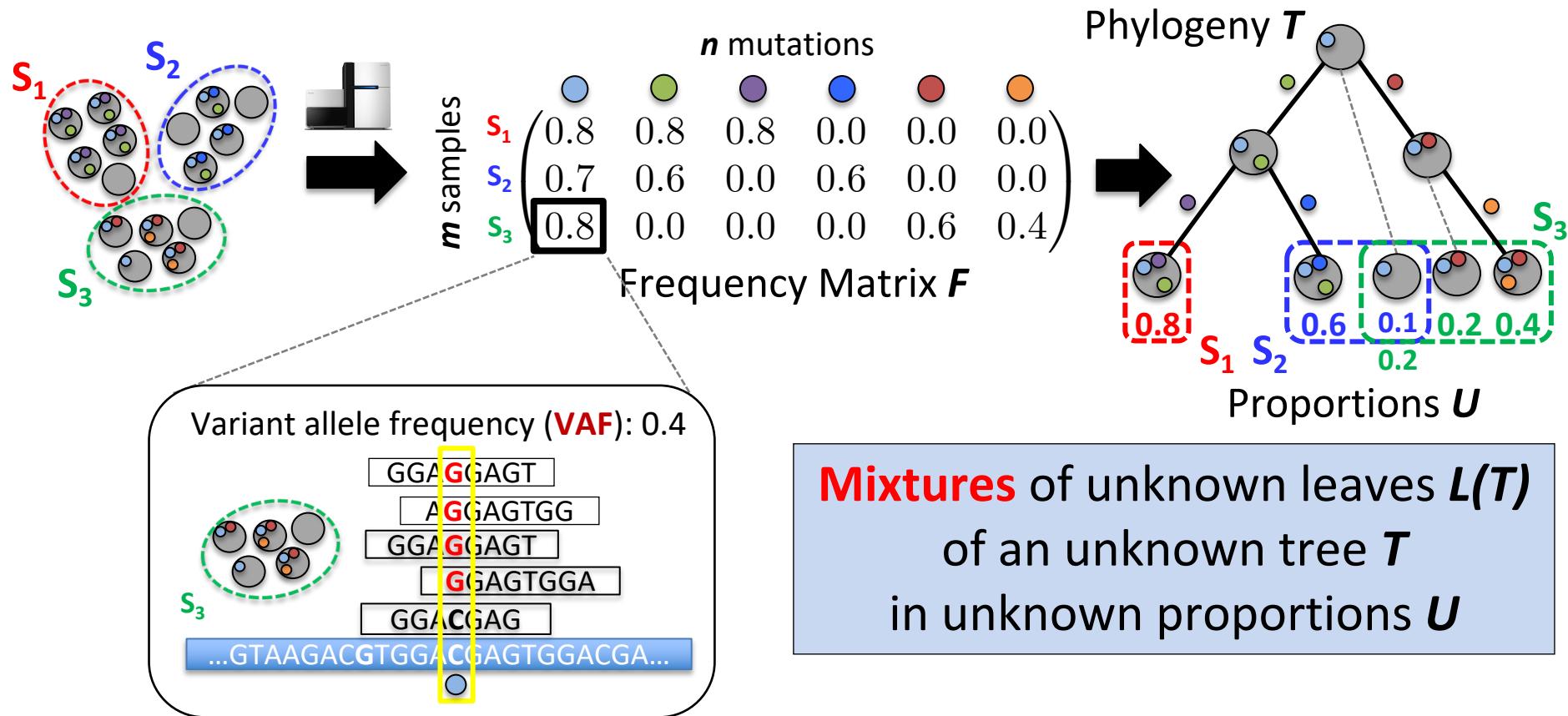
3. Summarizing solution space: [ISMB 2019]

- Multiple consensus tree problem

Sequencing and Tumor Phylogeny Inference



Sequencing and Tumor Phylogeny Inference



Tumor Phylogeny Inference: Given frequencies F , find phylogeny T and proportions U

Perfect Phylogeny Mixture

Assumptions:

- Infinite sites assumption:
a character changes state once
- Error-free data

$$\begin{array}{c}
 \text{Frequency Matrix } \mathbf{F} \\
 \begin{array}{c}
 \text{m samples} \\
 \text{n mutations} \\
 \begin{matrix} \textcolor{red}{S_1} & \begin{pmatrix} 0.8 & 0.8 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.7 & 0.6 & 0.0 & 0.6 & 0.0 & 0.0 \\ \textcolor{green}{S_3} & 0.8 & 0.0 & 0.0 & 0.0 & 0.6 & 0.4 \end{pmatrix} \\
 \textcolor{blue}{S_2} \end{matrix}
 \end{array} = \begin{array}{c}
 \text{m samples} \\
 \text{clones} \\
 \begin{matrix} \textcolor{red}{S_1} & \begin{pmatrix} 0.0 & 0.0 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.1 & 0.0 & 0.0 & 0.6 & 0.0 & 0.0 \\ \textcolor{green}{S_3} & 0.2 & 0.0 & 0.0 & 0.0 & 0.2 & 0.4 \end{pmatrix} \\
 \textcolor{blue}{S_2} \end{matrix}
 \end{array} \quad \text{Mixture Matrix } \mathbf{U} \\
 \begin{array}{c}
 \text{n mutations} \\
 \text{clones} \\
 \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \\
 \text{clones} \\
 \begin{matrix} \textcolor{red}{S_1} & \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \\
 \textcolor{blue}{S_2} \end{matrix}
 \end{array} \\
 \text{Restricted PP Matrix } \mathbf{B}
 \end{array}$$

Rows of \mathbf{U} are proportions:

$$u_{pj} \geq 0 \text{ and } \sum_j u_{pj} \leq 1$$

Perfect Phylogeny Theorem
 [Estabrook, 1971]
 [Gusfield, 1991]

Perfect Phylogeny Mixture: [El-Kebir*, Oesper* et al., 2015]
 Given \mathbf{F} , find \mathbf{U} and \mathbf{B} such that $\mathbf{F} = \mathbf{U} \mathbf{B}$

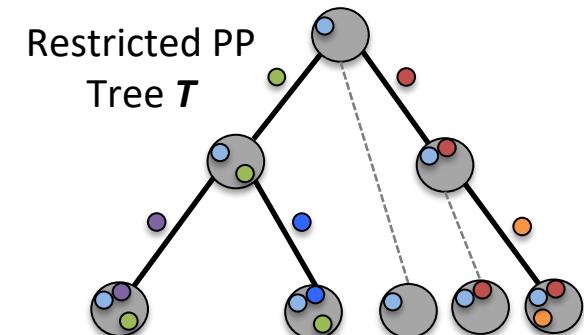
Previous Work

Variant of PPM:

TrAp [Strino *et al.*, 2013], PhyloSub [Jiao *et al.*, 2014]
 CITUP [Malikic *et al.*, 2015], BitPhylogeny [Yuan *et al.*, 2015]
 LICHHeE [Popic *et al.*, 2015], ...

$$\begin{matrix} m \text{ samples} \\ \text{Frequency Matrix } \mathbf{F} \end{matrix} = \begin{matrix} n \text{ mutations} \\ \mathbf{S}_1 \\ \mathbf{S}_2 \\ \mathbf{S}_3 \end{matrix} \left(\begin{matrix} 0.8 & 0.8 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.7 & 0.6 & 0.0 & 0.6 & 0.0 & 0.0 \\ 0.8 & 0.0 & 0.0 & 0.0 & 0.6 & 0.4 \end{matrix} \right)$$

$$\begin{matrix} m \text{ samples} \\ \text{clones} \\ \text{Mixture Matrix } \mathbf{U} \end{matrix} = \begin{matrix} n \text{ mutations} \\ \mathbf{S}_1 \\ \mathbf{S}_2 \\ \mathbf{S}_3 \end{matrix} \left(\begin{matrix} 0.0 & 0.0 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.1 & 0.0 & 0.0 & 0.6 & 0.0 & 0.0 \\ 0.2 & 0.0 & 0.0 & 0.0 & 0.2 & 0.4 \end{matrix} \right) \quad \begin{matrix} n \text{ mutations} \\ \text{clones} \\ \text{Restricted PP Matrix } \mathbf{B} \end{matrix}$$



1-1 Equivalent

Rows of \mathbf{U} are proportions:

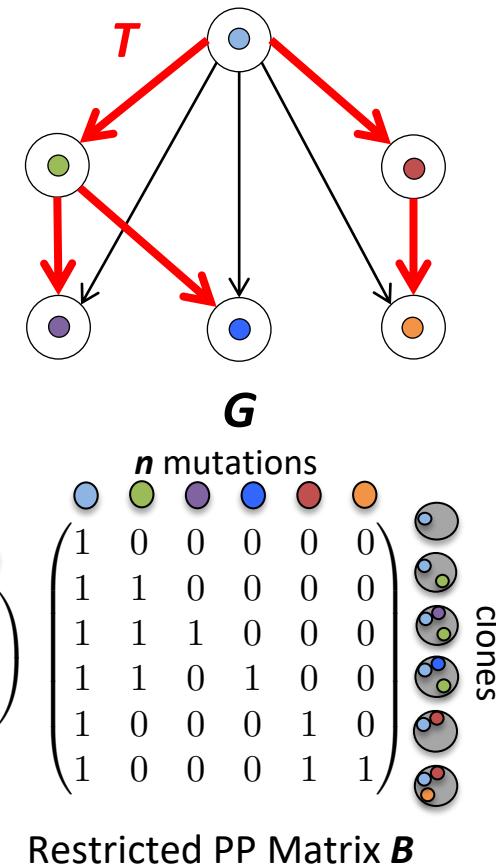
$$u_{pj} \geq 0 \text{ and } \sum_j u_{pj} \leq 1$$

Perfect Phylogeny Theorem
 [Estabrook, 1971]
 [Gusfield, 1991]

Perfect Phylogeny Mixture: [El-Kebir*, Oesper* et al., 2015]
 Given \mathbf{F} , find \mathbf{U} and \mathbf{B} such that $\mathbf{F} = \mathbf{U} \mathbf{B}$

Combinatorial Characterization

- Frequency $f_{p,i}$ is mass of subtree rooted at node that introduced i
- Usage $u_{p,i}$ is mass of node that introduced i



$$\begin{matrix} & \text{n mutations} \\ m \text{ samples} & \begin{matrix} S_1 & S_2 & S_3 \end{matrix} \end{matrix} = \begin{matrix} & \text{n mutations} \\ m \text{ samples} & \begin{matrix} S_1 & S_2 & S_3 \end{matrix} \end{matrix}$$

$$\text{Frequency Matrix } \mathbf{F}$$

$$\begin{pmatrix} 0.8 & 0.8 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.7 & 0.6 & 0.0 & 0.6 & 0.0 & 0.0 \\ 0.8 & 0.0 & 0.0 & 0.0 & 0.6 & 0.4 \end{pmatrix}$$

$$\text{Mixture Matrix } \mathbf{U}$$

$$\begin{pmatrix} 0.0 & 0.0 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.1 & 0.0 & 0.0 & 0.6 & 0.0 & 0.0 \\ 0.2 & 0.0 & 0.0 & 0.0 & 0.2 & 0.4 \end{pmatrix}$$

$$\text{Restricted PP Matrix } \mathbf{B}$$

$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

Theorem 1:

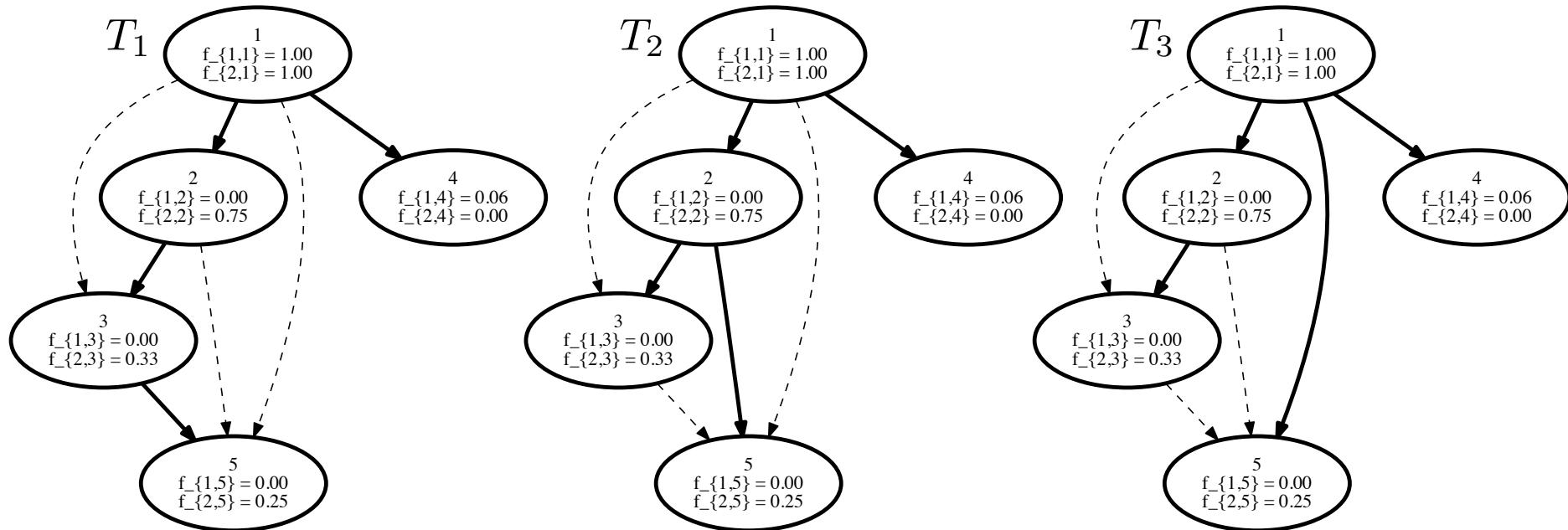
T is a solution to the PPM if and only if T is a spanning tree of G satisfying the sum condition

Theorem 2:

PPM is NP-complete even for $m=2$

Perfect Phylogeny Mixture: [El-Kebir*, Oesper* et al., 2015]
Given \mathbf{F} , find \mathbf{U} and \mathbf{B} such that $\mathbf{F} = \mathbf{U}\mathbf{B}$

Non-uniqueness of Solutions to PPM



$$F = \begin{pmatrix} 1 & 0 & 0 & 0.06 & 0 \\ 1 & 0.75 & 0.33 & 0 & 0.25 \end{pmatrix}$$

Question 1: Can we determine the number of solutions?

Question 2: Can sample solutions uniformly at random?
21

On the Complexity of #PPM (new results)

Question 1: Can we determine the number of solutions?

Question 2: Can sample solutions uniformly at random?

#PPM: Given F , count the number of pairs (U, B) composed of mixture matrix U and perfect phylogeny matrix B such that $F = UB$

On the Complexity of #PPM (new results)

Question 1: Can we determine the number of solutions?

Question 2: Can sample solutions uniformly at random?

#PPM: Given F , count the number of pairs (U, B) composed of mixture matrix U and perfect phylogeny matrix B such that $F = UB$

#P is the complexity class of counting problems whose decision problems are in NP

Every problem in #P can be reduced in polynomial time to any problem in #P-complete, preserving cardinalities

On the Complexity of #PPM (new results)

Question 1: Can we determine the number of solutions?

Question 2: Can sample solutions uniformly at random?

#PPM: Given F , count the number of pairs (U, B) composed of mixture matrix U and perfect phylogeny matrix B such that $F = UB$

#P is the complexity class of counting problems whose decision problems are in NP

Every problem in #P can be reduced in polynomial time to any problem in #P-complete, preserving cardinalities

Theorem: #PPM is #P-complete

Theorem: There is no FPRAS for #PPM

Theorem: There is no FPAUS for PPM



Yuanyuan Qi

Outline

1. Background and theory: [RECOMB-CG 2018]

- Perfect Phylogeny Mixture (PPM) problem
- Combinatorial characterization of solutions
- #PPM: exact counting and uniform sampling

2. Simulation results: [RECOMB-CG 2018]

- What contributes to non-uniqueness?
- How to reduce non-uniqueness?
- How does non-uniqueness affect current methods?

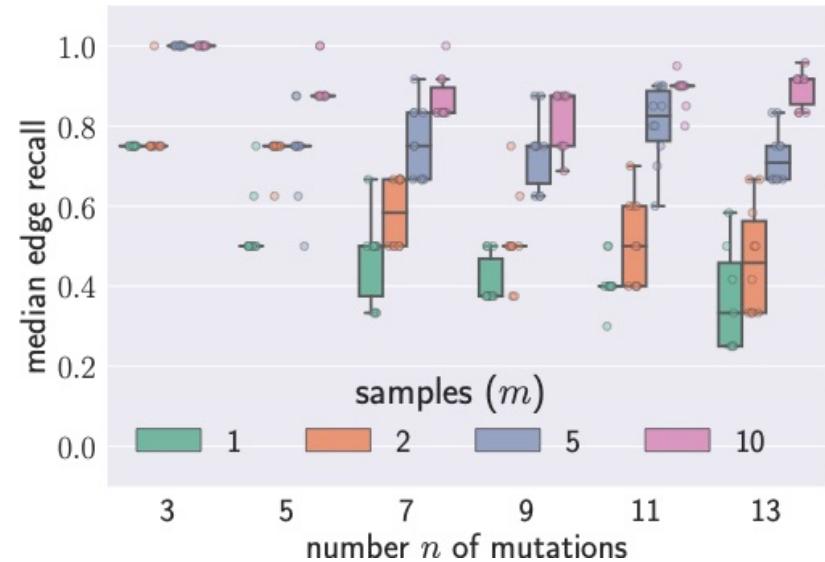
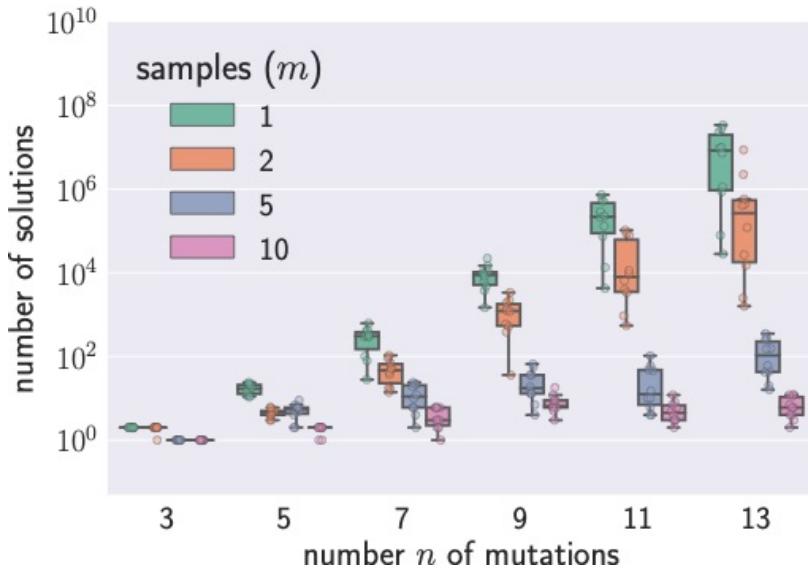


Dikshant Pradhan

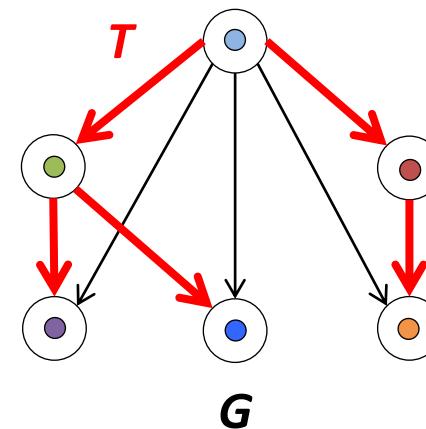
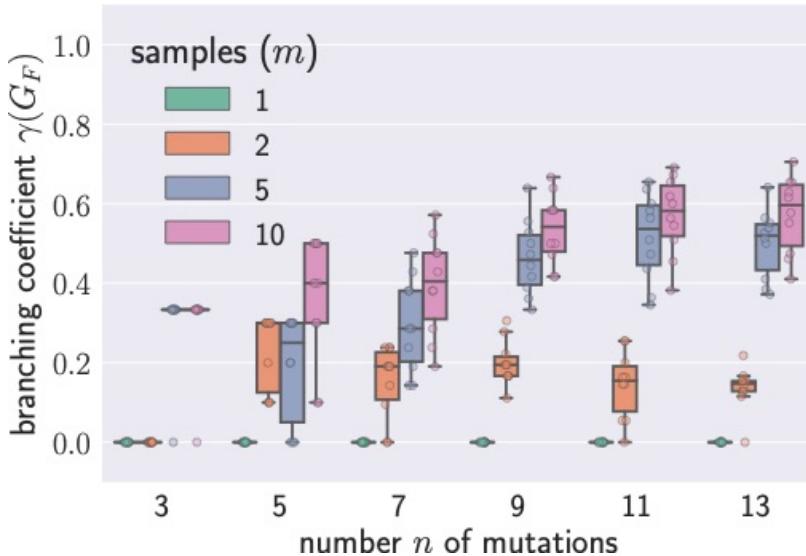
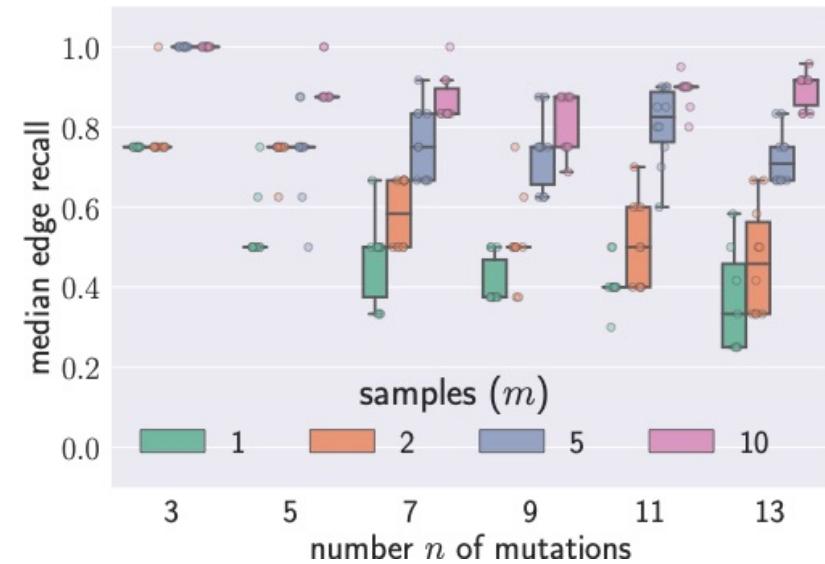
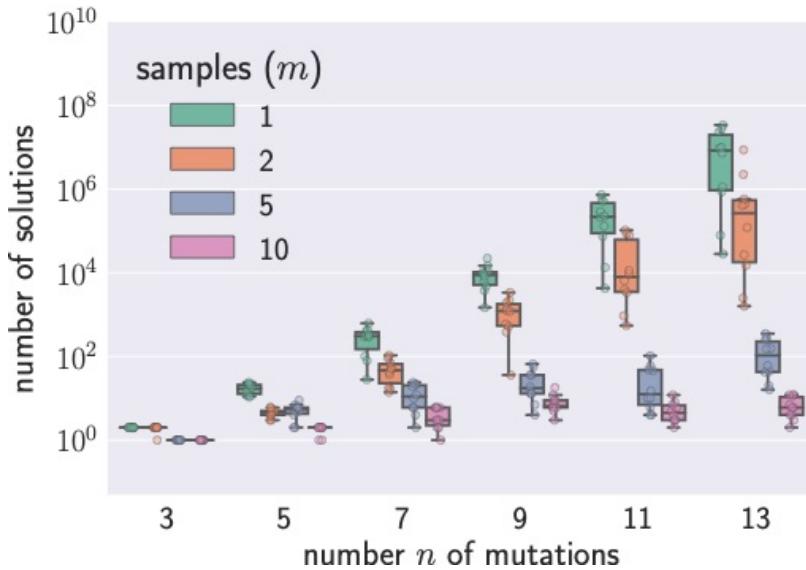
3. Summarizing solution space: [ISMB 2019]

- Multiple consensus tree problem

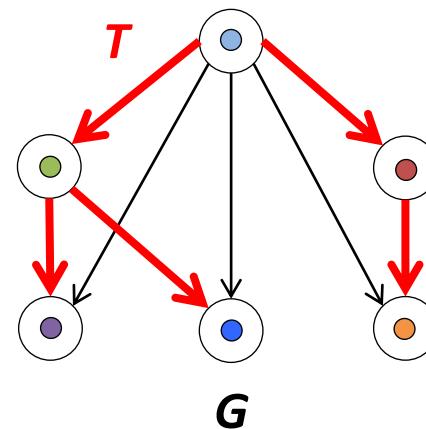
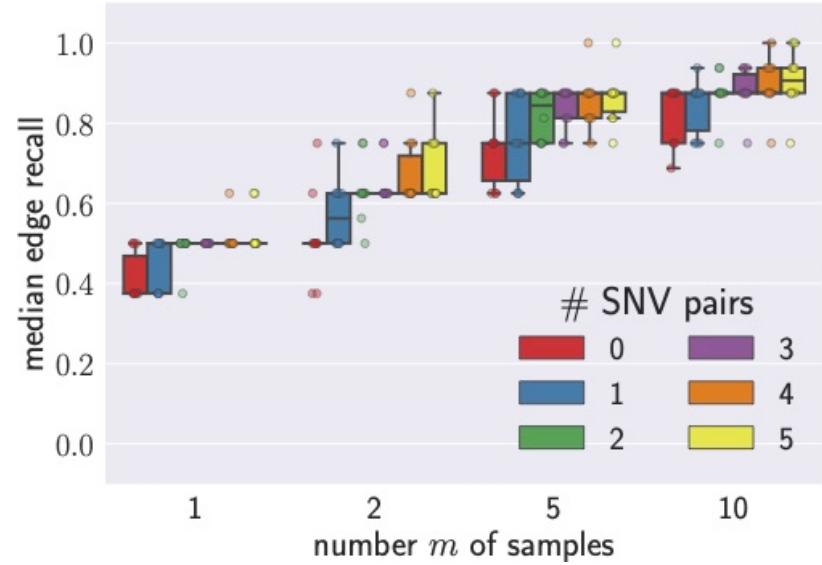
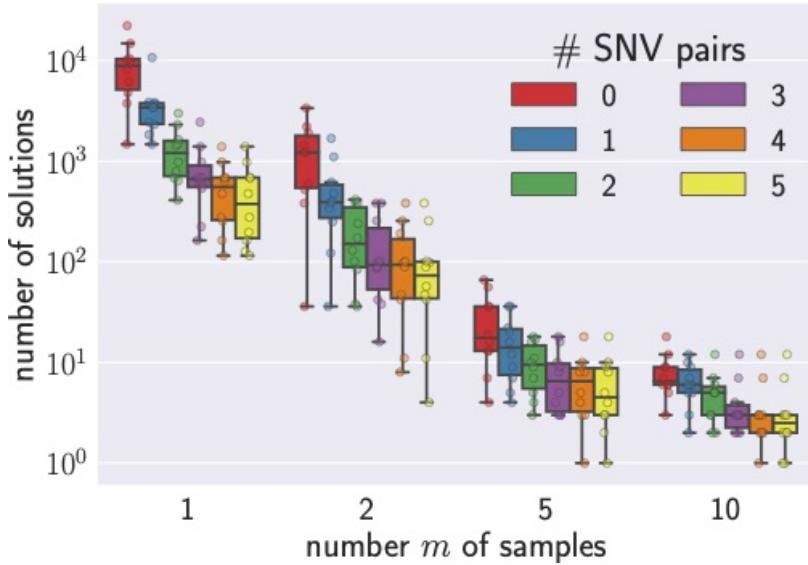
What Contributors to Non-uniqueness?



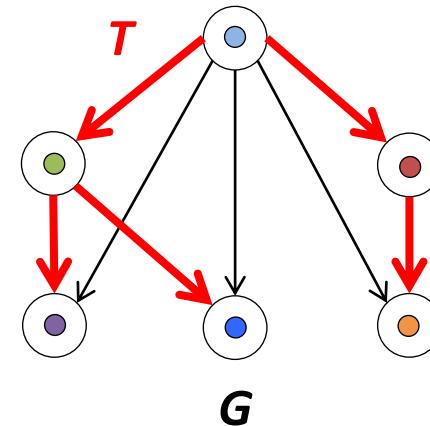
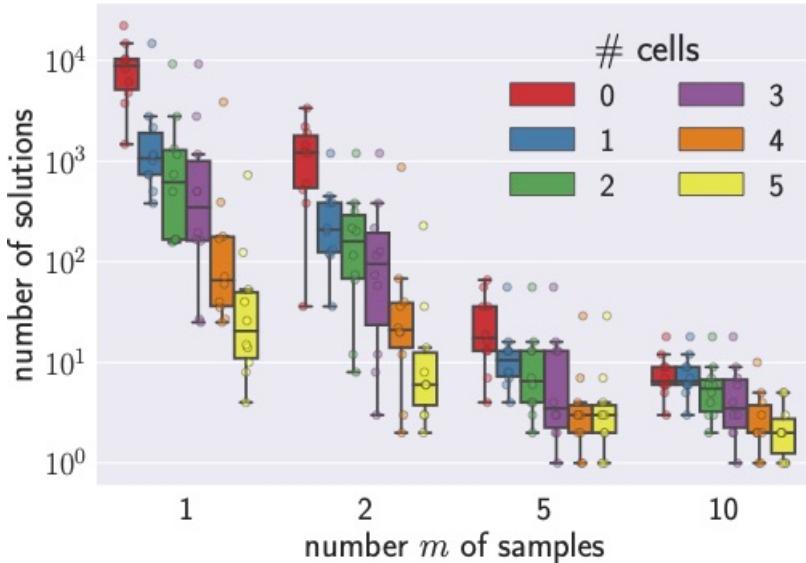
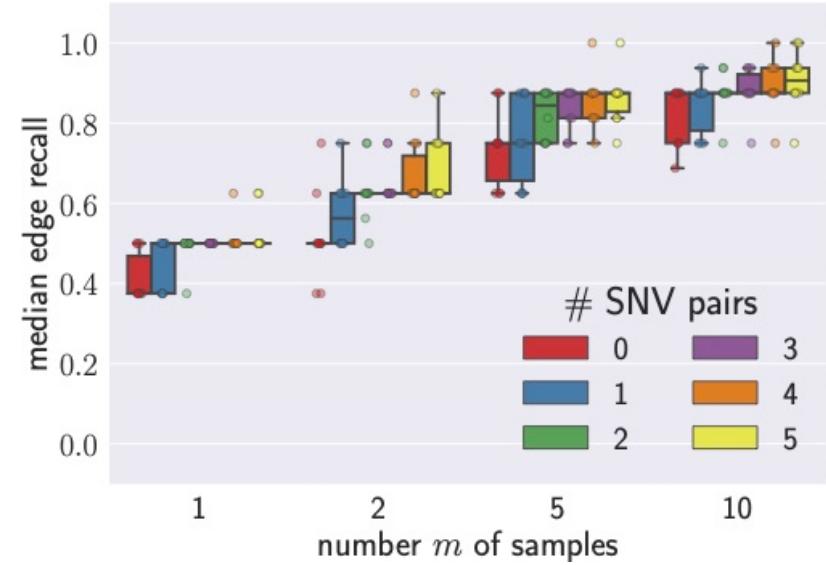
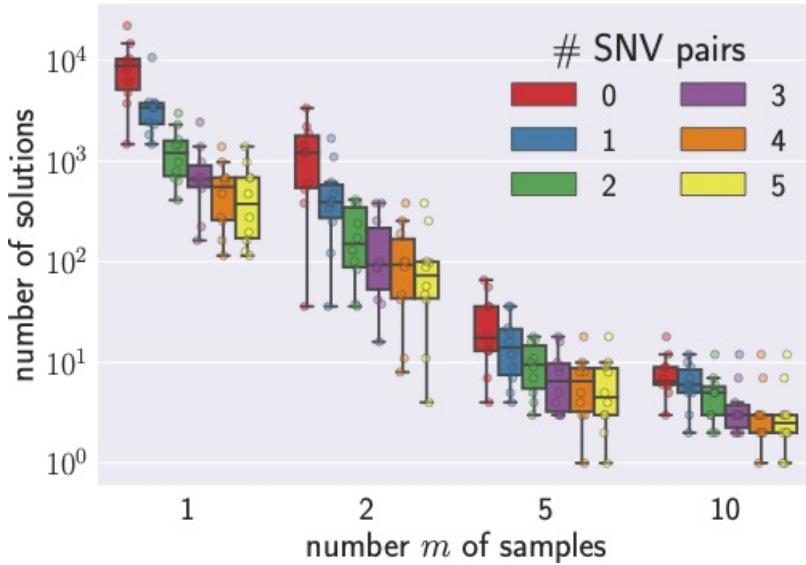
What Contributors to Non-uniqueness?



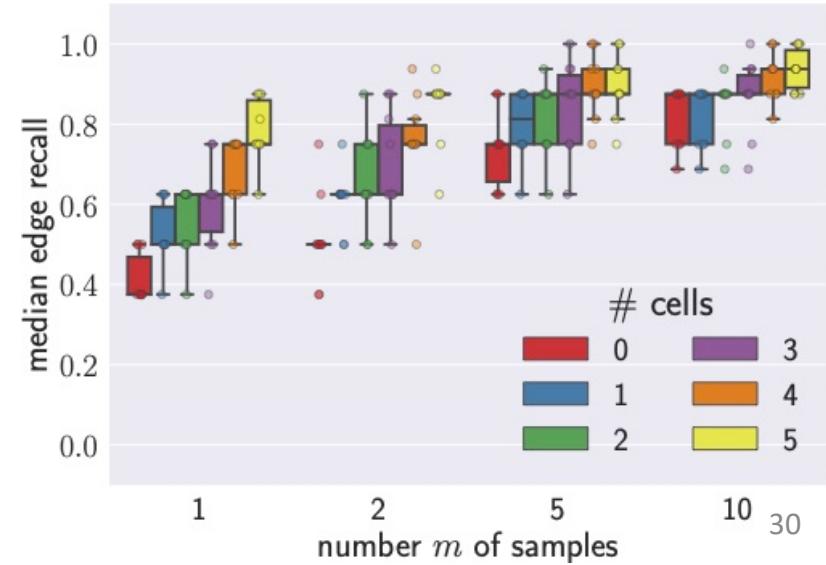
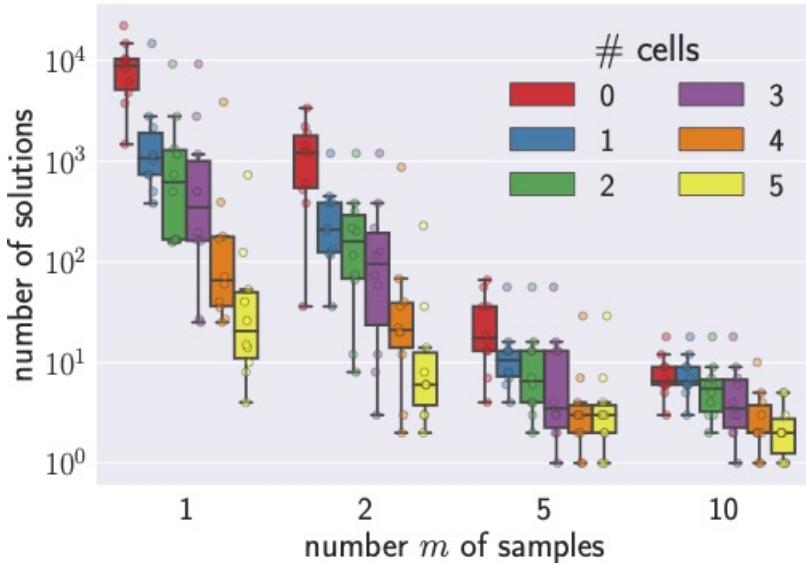
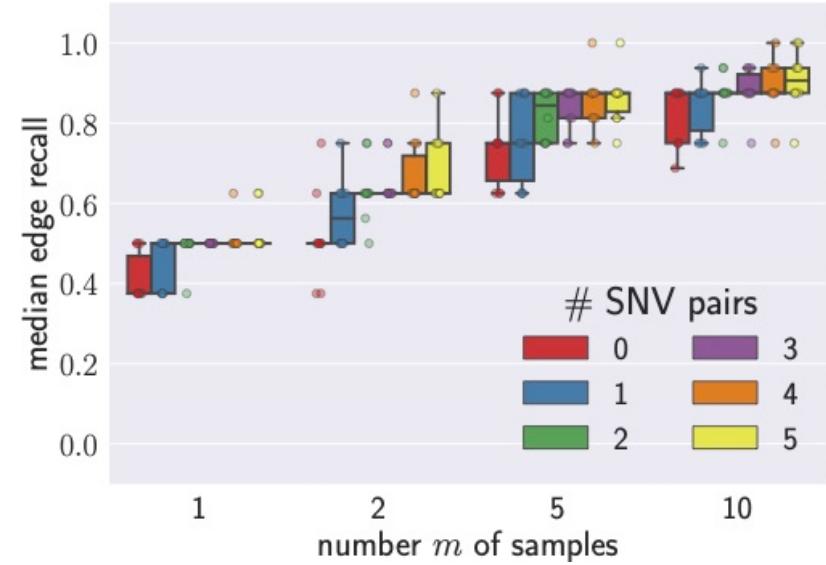
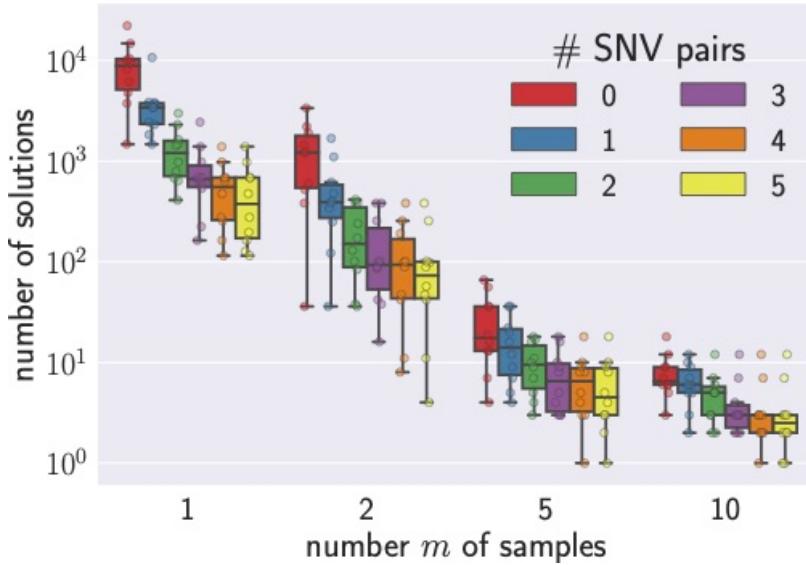
How to Reduce Non-Uniqueness?



How to Reduce Non-Uniqueness?



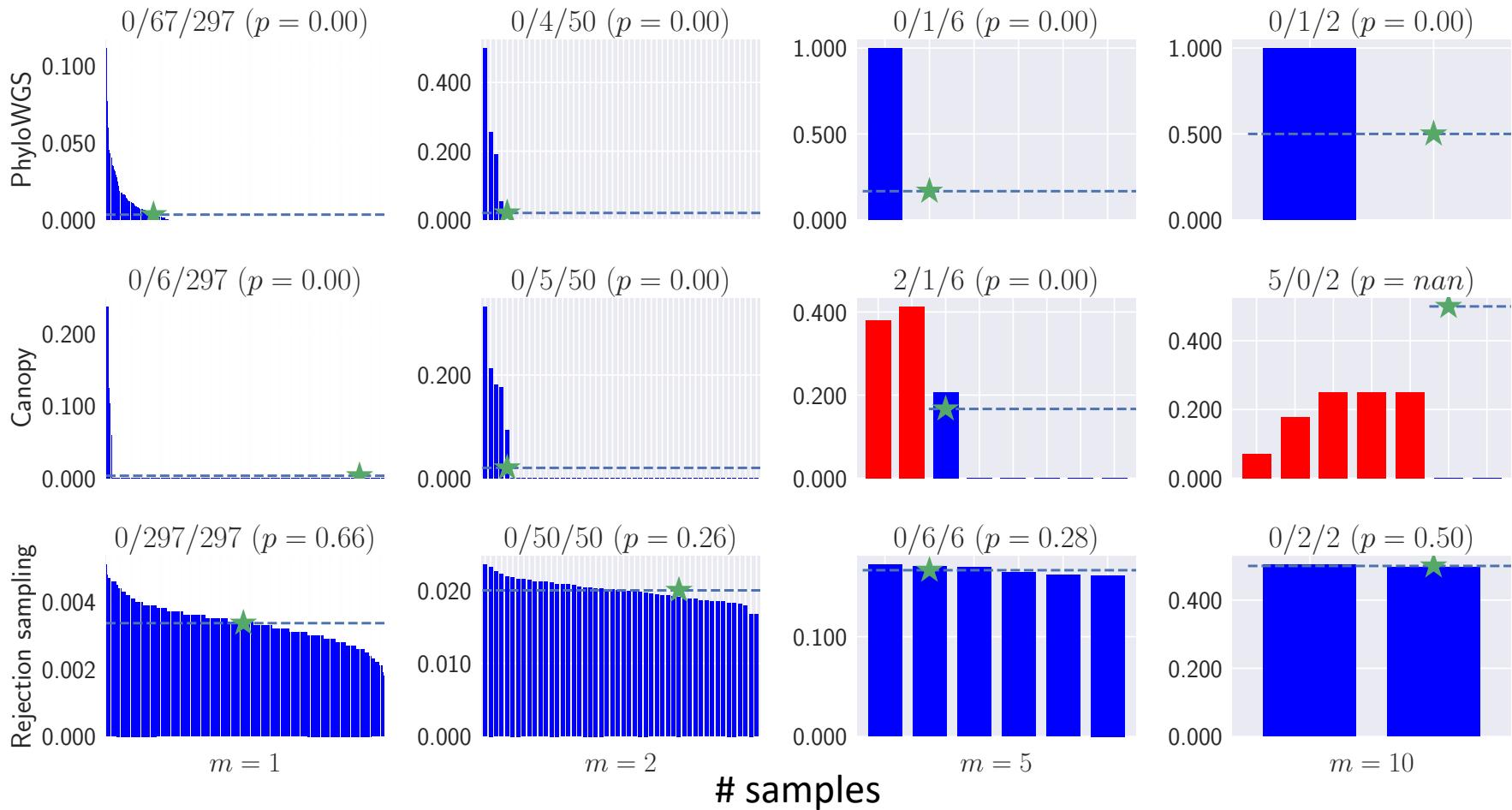
How to Reduce Non-Uniqueness?



How Does Non-uniqueness affect Methods?

Two current MCMC methods using default parameters:

- PhyloWGS, Deshwar et al., Genom. Biol., 2015 [10,000 samples]
- Canopy, Jiang et al., PNAS, 2016 [~300 samples]



Outline

1. Background and theory: [RECOMB-CG 2018]

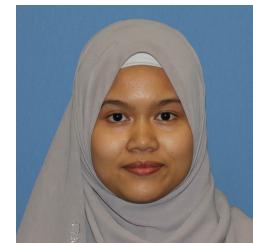
- Perfect Phylogeny Mixture (PPM) problem
- Combinatorial characterization of solutions
- #PPM: exact counting and uniform sampling

2. Simulation results: [RECOMB-CG 2018]

- What contributes to non-uniqueness?
- How to reduce non-uniqueness?
- How does non-uniqueness affect current methods?

3. Summarizing solution space: [ISMB 2019]

- Multiple consensus tree problem



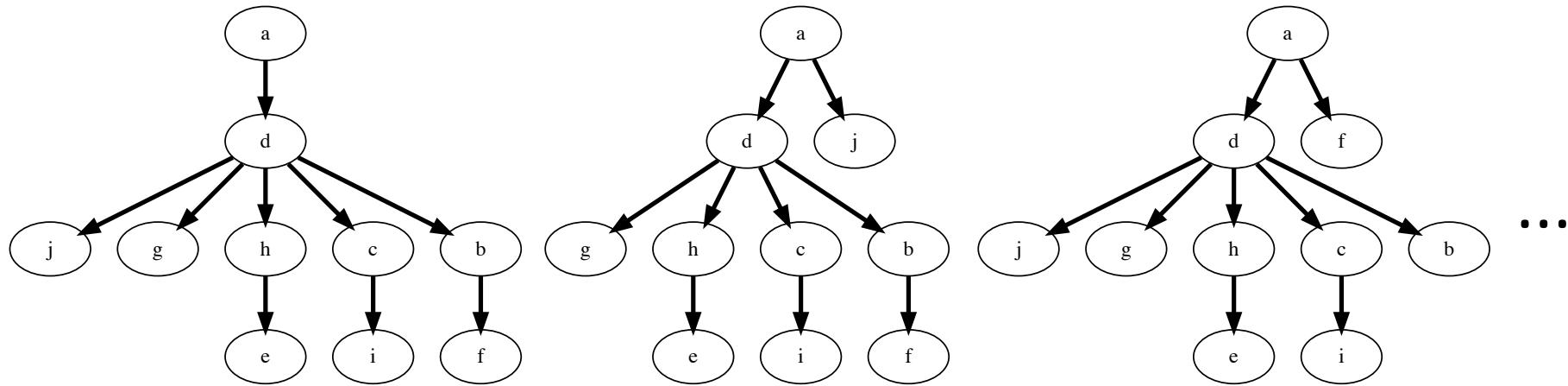
Nuraini Aguse



Yuanyuan Qi

Lung Cancer Patient: CRUK0037

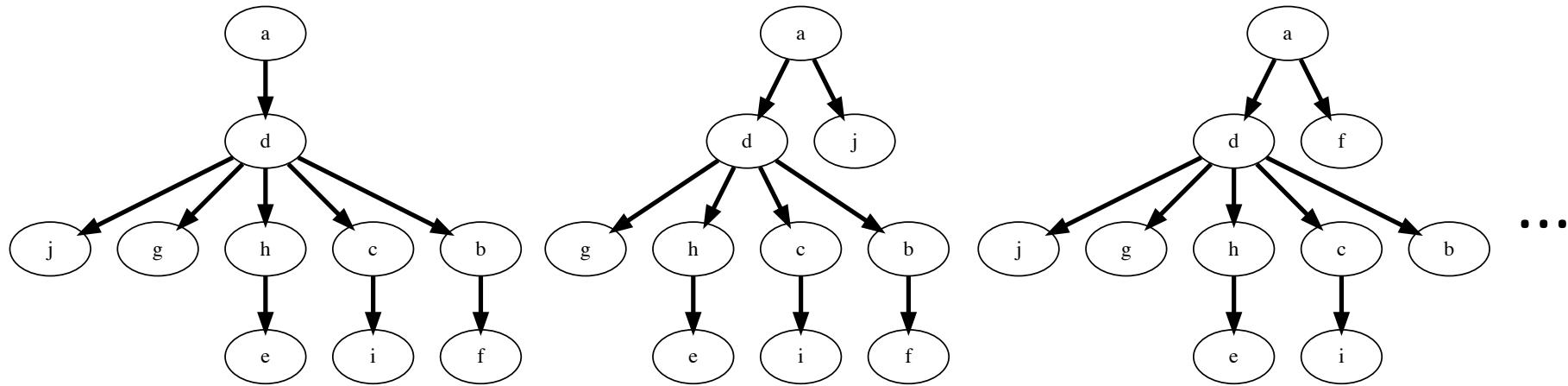
Jamal-Hanjani et al. (2017). *New England Journal of Medicine*, 376(22), 2109–2121.



Authors inferred 17 trees

Lung Cancer Patient: CRUK0037

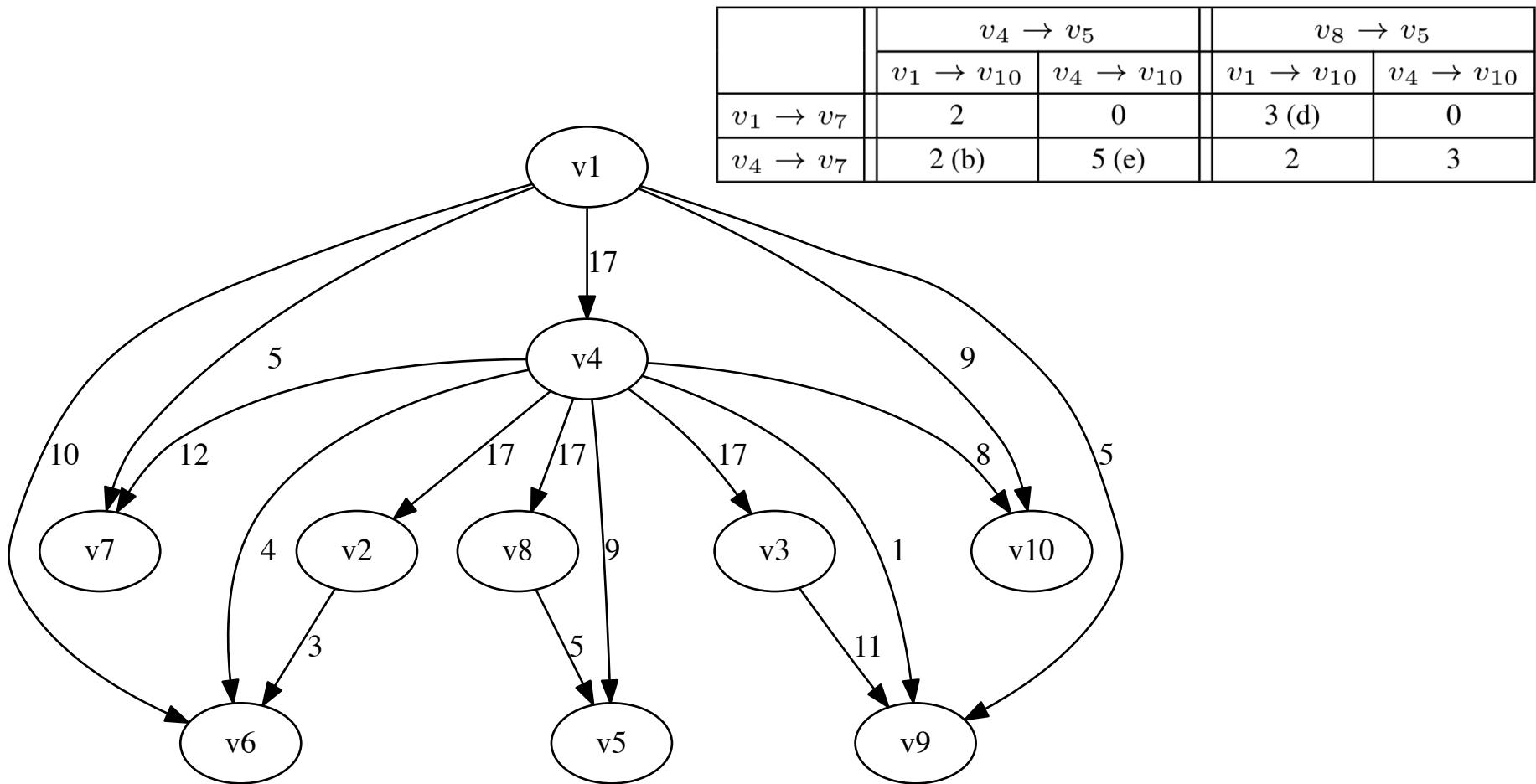
Jamal-Hanjani et al. (2017). *New England Journal of Medicine*, 376(22), 2109–2121.



Authors inferred 17 trees

Question: How to summarize solution space in order to remove inference errors and identify dependencies among mutations?

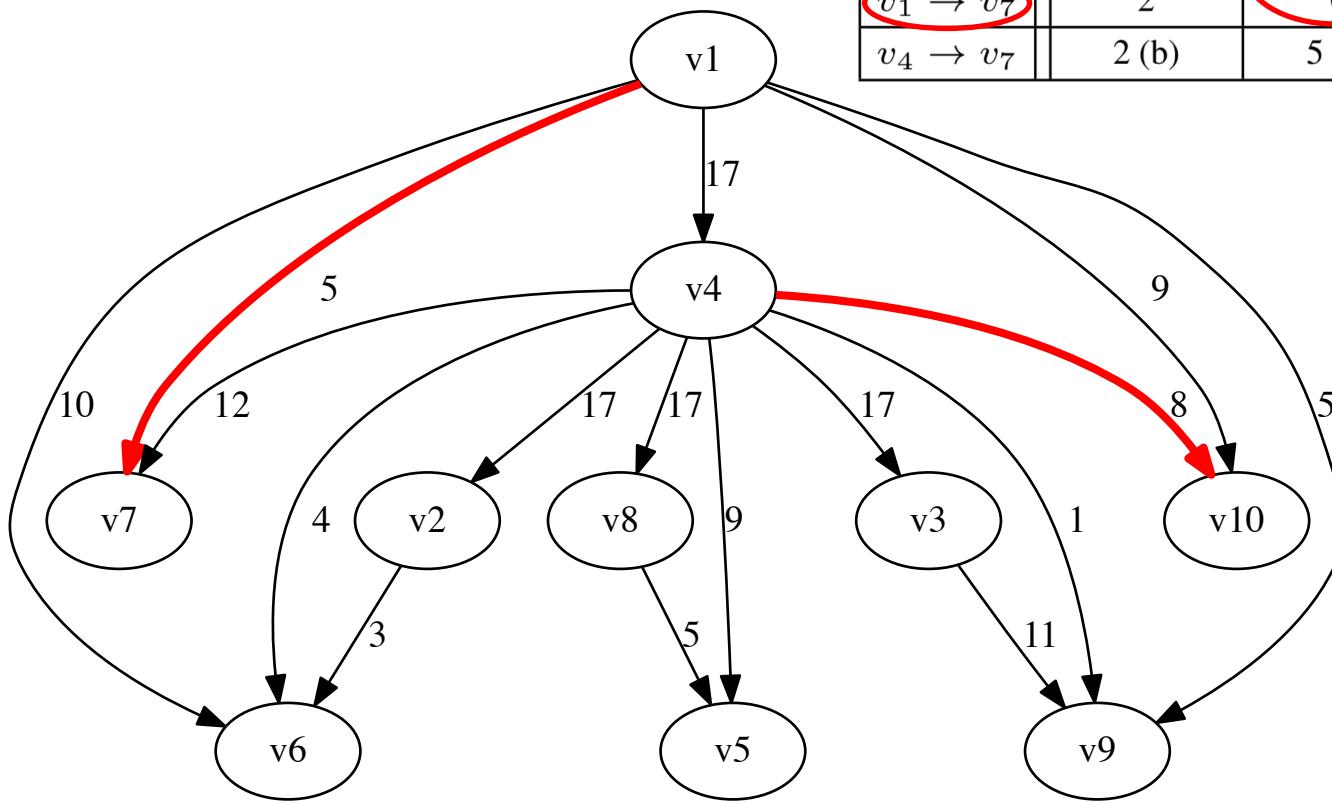
Parent-child Graph: Union of all Edges



	$v_4 \rightarrow v_5$		$v_8 \rightarrow v_5$	
	$v_1 \rightarrow v_{10}$	$v_4 \rightarrow v_{10}$	$v_1 \rightarrow v_{10}$	$v_4 \rightarrow v_{10}$
$v_1 \rightarrow v_7$	2	0	3 (d)	0
$v_4 \rightarrow v_7$	2 (b)	5 (e)	2	3

Parent-child Graph: Union of all Edges

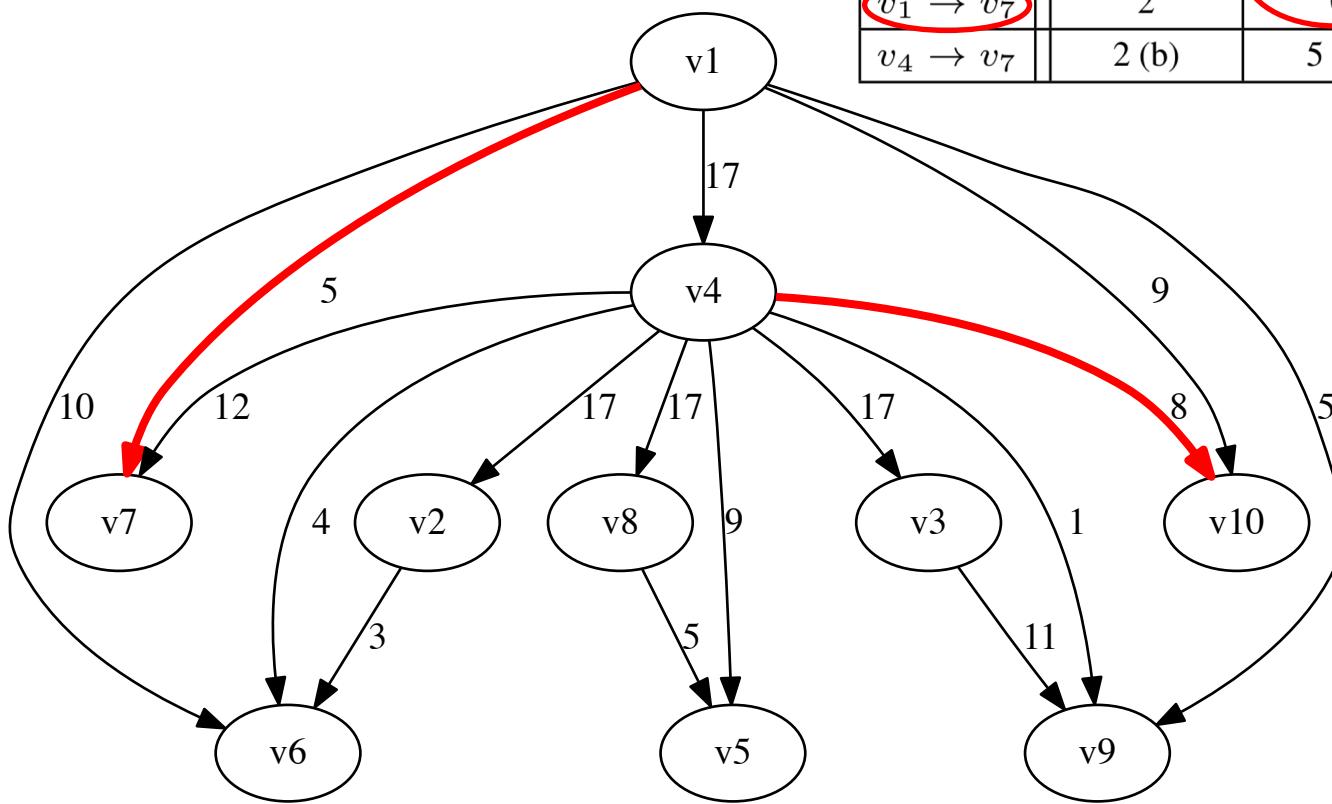
	$v_4 \rightarrow v_5$		$v_8 \rightarrow v_5$	
	$v_1 \rightarrow v_{10}$	$v_4 \rightarrow v_{10}$	$v_1 \rightarrow v_{10}$	$v_4 \rightarrow v_{10}$
$v_1 \rightarrow v_7$	2	0	3 (d)	0
$v_4 \rightarrow v_7$	2 (b)	5 (e)	2	3



The parent-child graph does capture patterns of mutual exclusivity

Parent-child Graph: Union of all Edges

	$v_4 \rightarrow v_5$		$v_8 \rightarrow v_5$	
	$v_1 \rightarrow v_{10}$	$v_4 \rightarrow v_{10}$	$v_1 \rightarrow v_{10}$	$v_4 \rightarrow v_{10}$
$v_1 \rightarrow v_7$	2	0	3 (d)	0
$v_4 \rightarrow v_7$	2 (b)	5 (e)	2	3

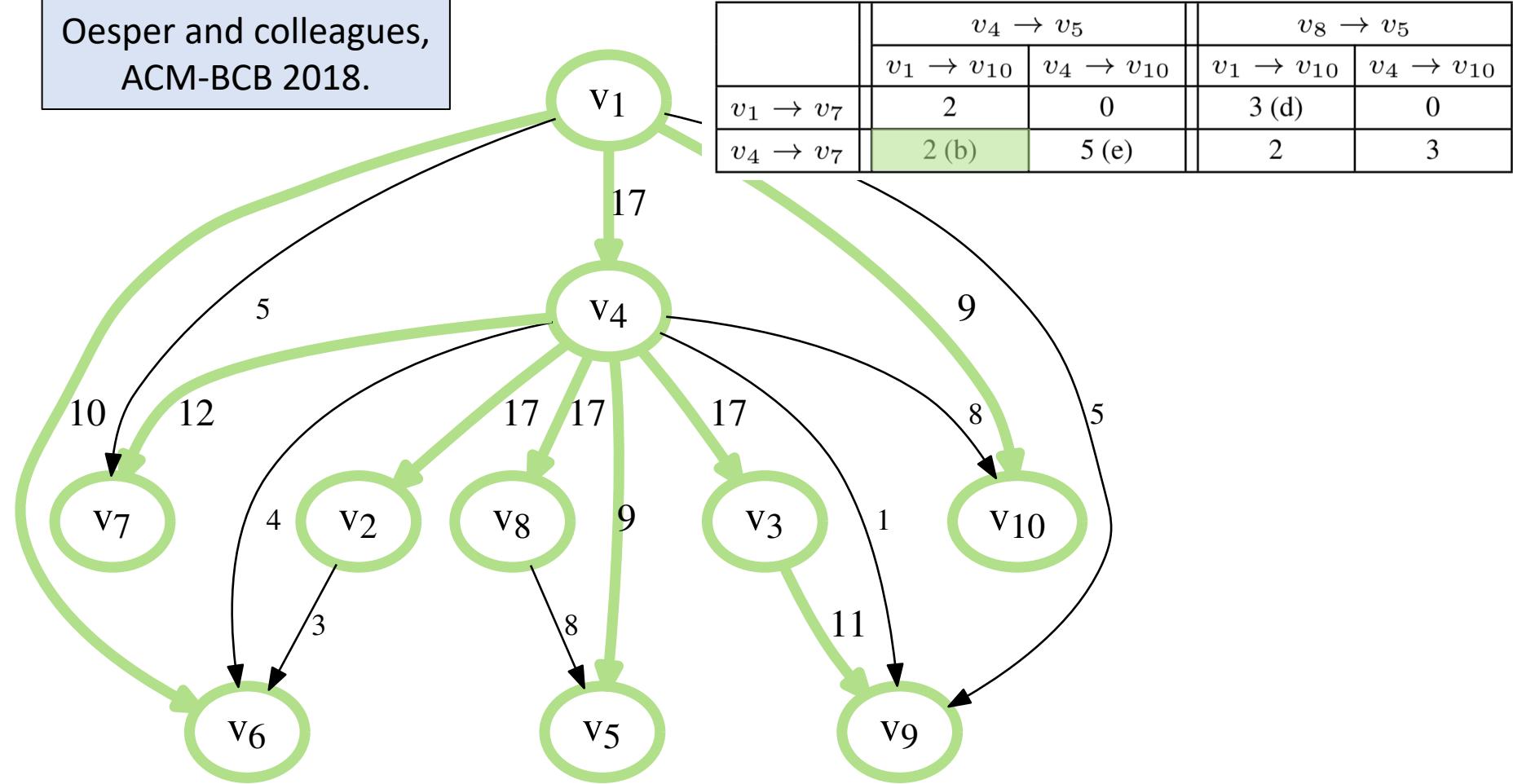


The parent-child graph does capture patterns of mutual exclusivity

Question: Can we infer a single consensus tree?

Single Consensus Tree: Max Weight Spanning Tree

Oesper and colleagues,
ACM-BCB 2018.

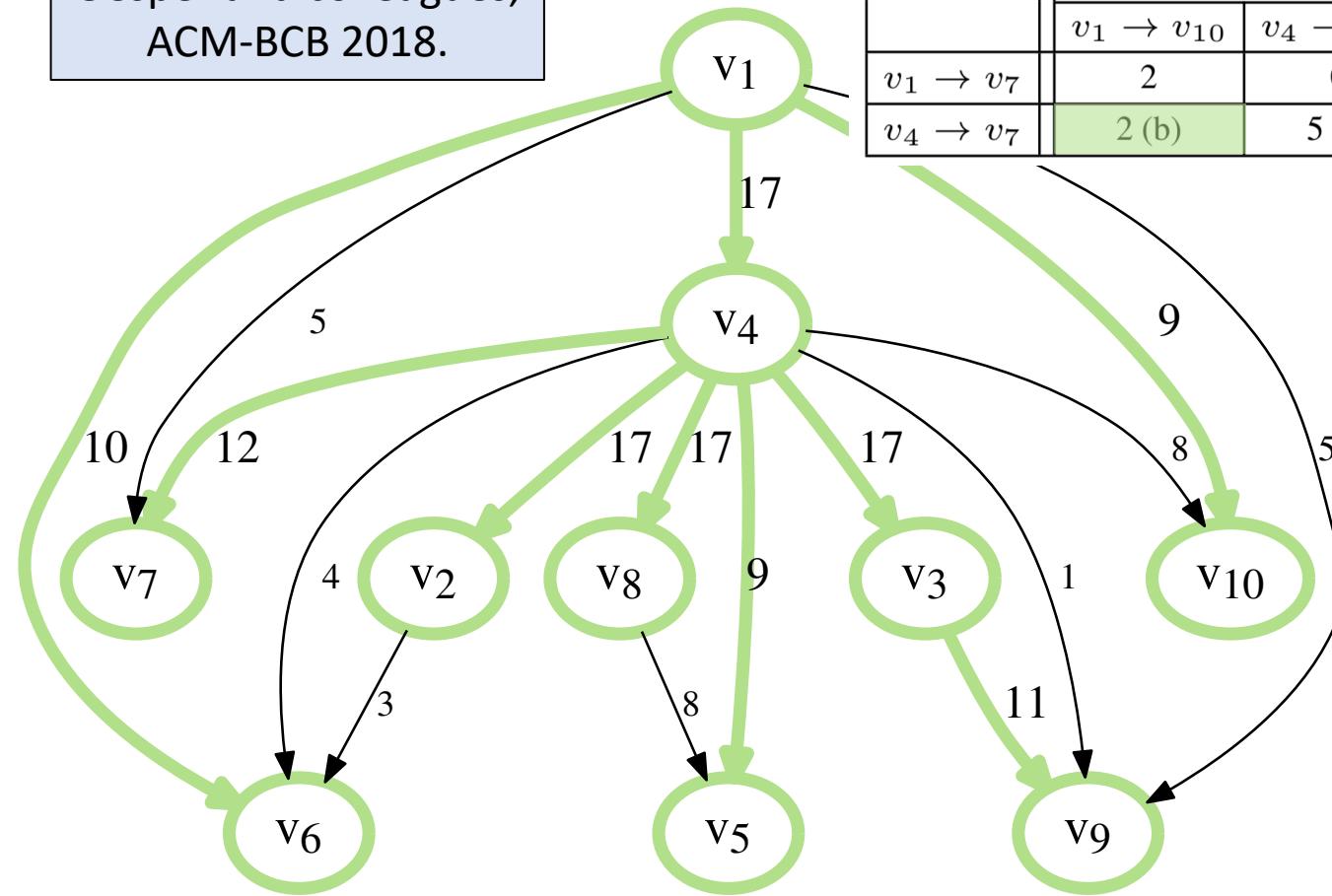


	$v_4 \rightarrow v_5$		$v_8 \rightarrow v_5$	
	$v_1 \rightarrow v_{10}$	$v_4 \rightarrow v_{10}$	$v_1 \rightarrow v_{10}$	$v_4 \rightarrow v_{10}$
$v_1 \rightarrow v_7$	2	0	3 (d)	0
$v_4 \rightarrow v_7$	2 (b)	5 (e)	2	3

Single Consensus Tree: Max Weight Spanning Tree

Oesper and colleagues,
ACM-BCB 2018.

	$v_4 \rightarrow v_5$		$v_8 \rightarrow v_5$	
	$v_1 \rightarrow v_{10}$	$v_4 \rightarrow v_{10}$	$v_1 \rightarrow v_{10}$	$v_4 \rightarrow v_{10}$
$v_1 \rightarrow v_7$	2	0	3 (d)	0
$v_4 \rightarrow v_7$	2 (b)	5 (e)	2	3

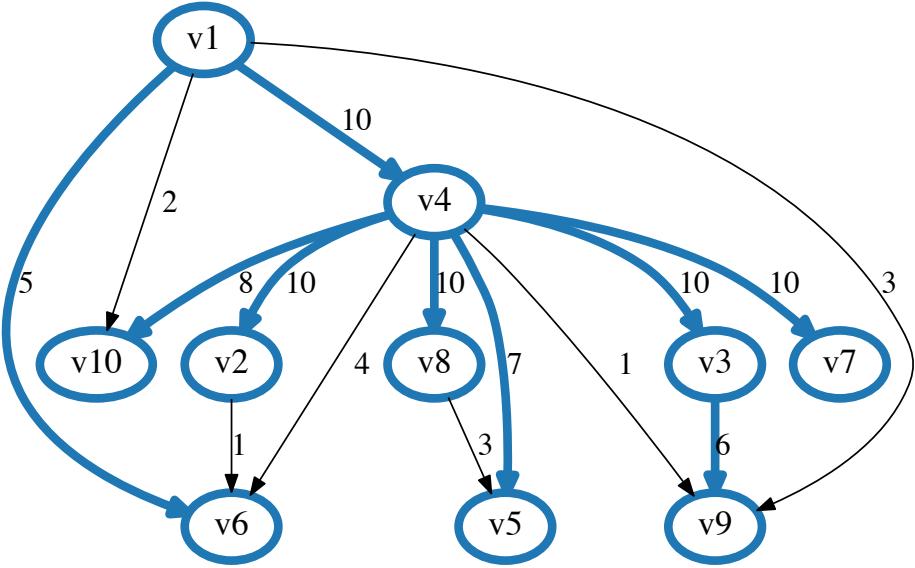
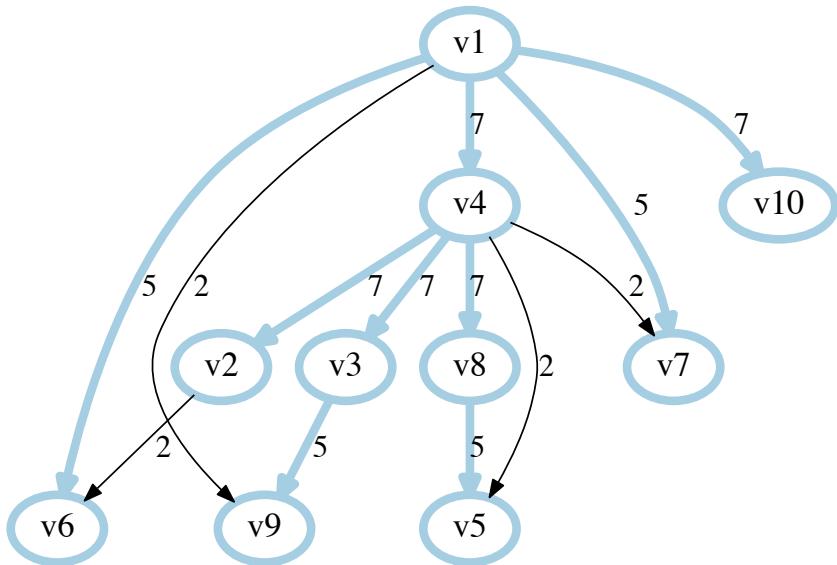


Inaccurate summary for diverse solution spaces

Question: How about inferring multiple consensus trees?

Multiple Consensus Trees

Simultaneous clustering and consensus tree inference

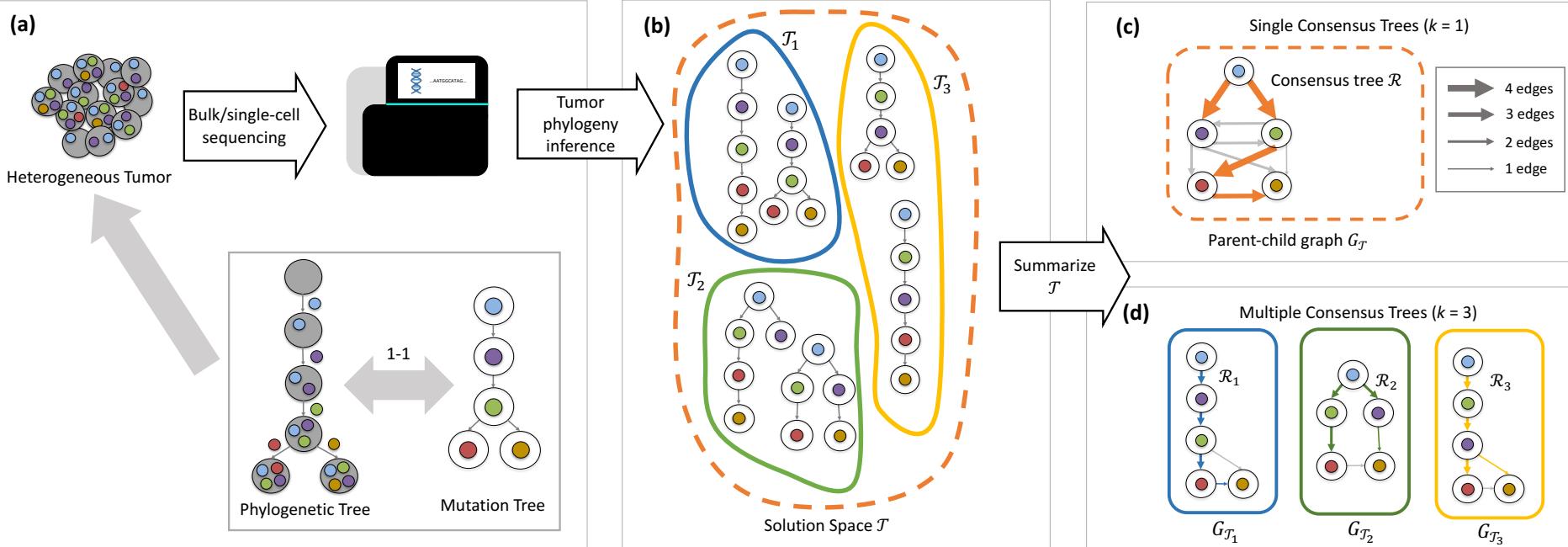


Multiple Consensus Trees (MCT): [ISMB 2019]

Given trees $\mathcal{T} = \{T_1, \dots, T_n\}$, find surjective clustering $\sigma : [n] \rightarrow [k]$ and consensus trees $\mathcal{R} = \{R_1, \dots, R_k\}$ such that $\sum_{i=1}^n d(T_i, R_{\sigma(i)})$ is minimum

Multiple Consensus Trees (MCT): [ISMB 2019]

Given trees $\mathcal{T} = \{T_1, \dots, T_n\}$, find surjective clustering $\sigma : [n] \rightarrow [k]$ and consensus trees $\mathcal{R} = \{R_1, \dots, R_k\}$ such that $\sum_{i=1}^n d(T_i, R_{\sigma(i)})$ is minimum



- Characterize combinatorial structure of optimal solutions
- Show that MCT is NP-hard for general k
- Introduce an MILP for solving the problem for small instance sizes
- Introduce a heuristic that returns optimal solution in most cases

Conclusion

1. Background and theory: [RECOMB-CG 2018]

- Perfect Phylogeny Mixture (PPM) problem
- Combinatorial characterization of solutions
- #PPM: exact counting and uniform sampling

2. Simulation results: [RECOMB-CG 2018]

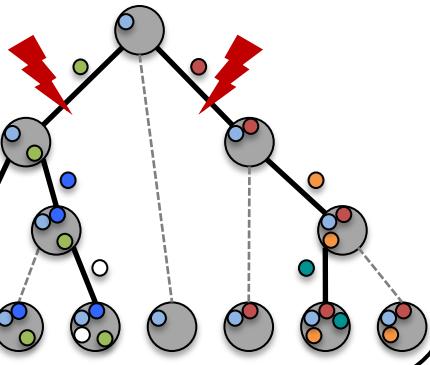
- What contributes to non-uniqueness?
- How to reduce non-uniqueness?
- How does non-uniqueness affect current methods?

3. Summarizing solution space: [ISMB 2019]

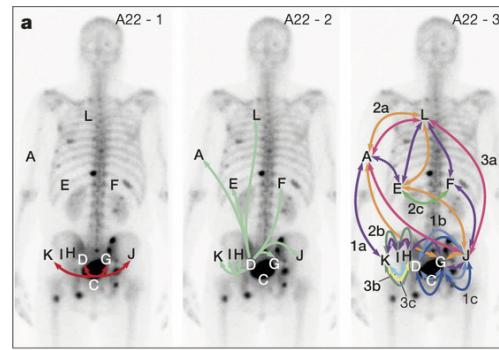
- Multiple consensus tree problem

Outlook

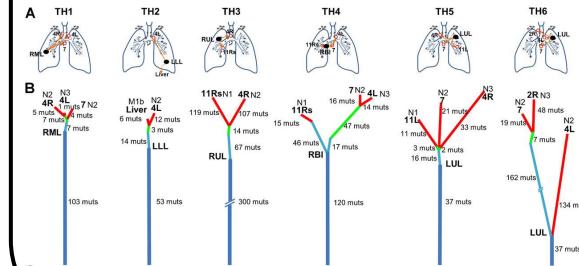
Identify targets for treatment



Understand metastatic development



Recognize common patterns of tumor evolution across patients



Downstream analyses in cancer genomics **critically rely** on accurate tumor phylogeny inference

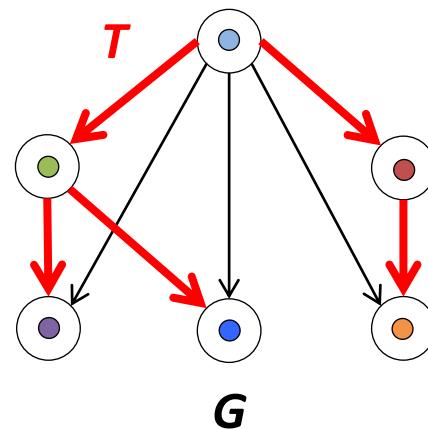
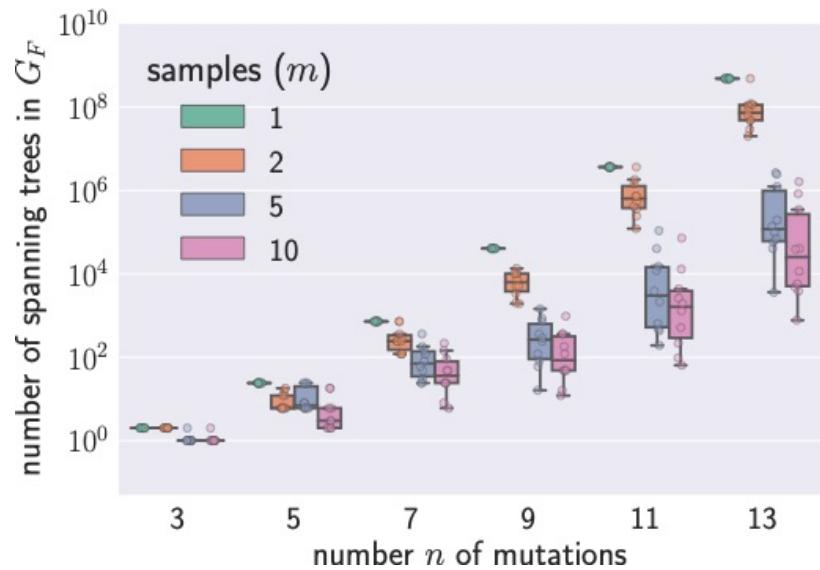
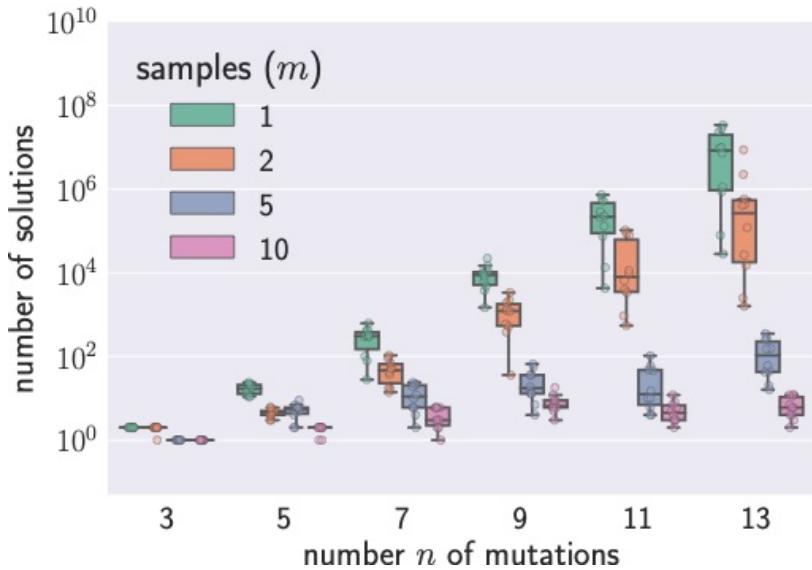
Challenge:

Novel algorithms that sample **uniformly at random** from the space of PPM solutions

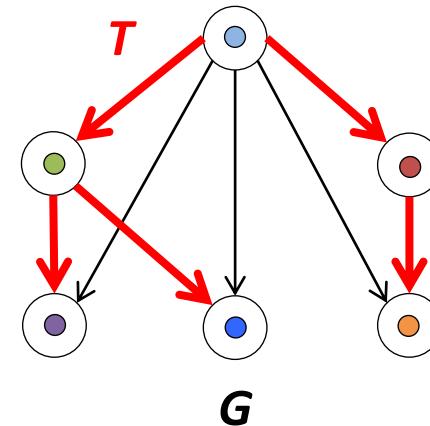
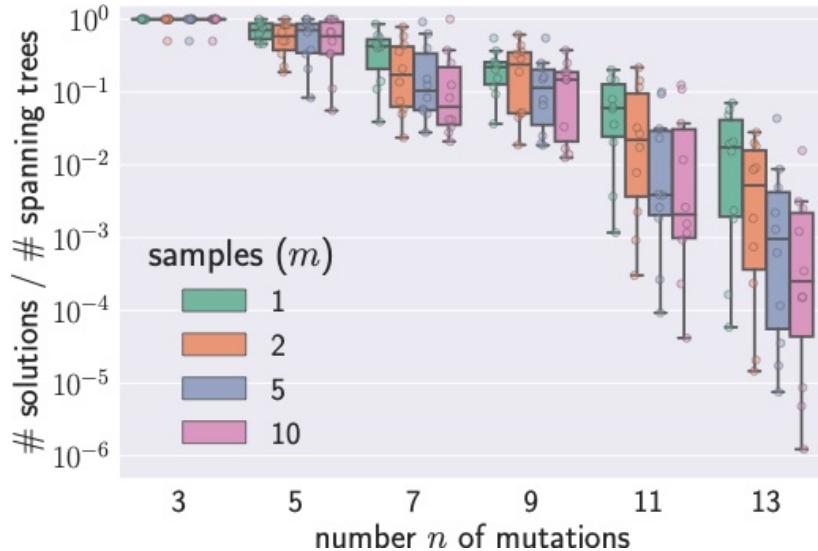
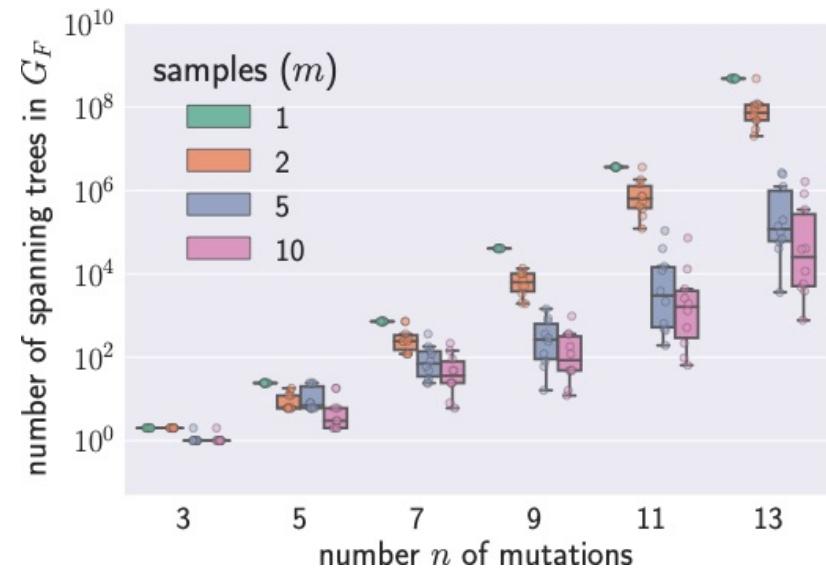
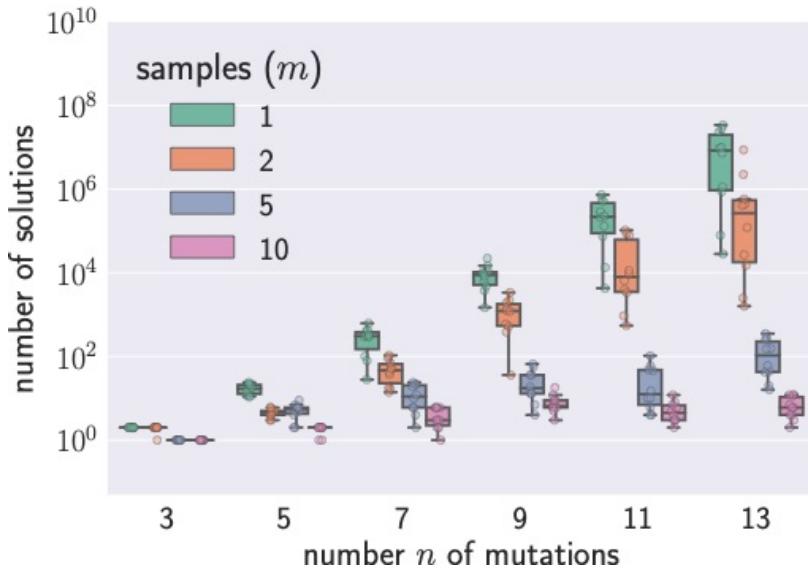
Acknowledgments

- Yuanyuan Qi
- Nuraini Aguse
- Dikshant Pradhan
- Experiments were run on NCSA's Blue Waters supercomputer
- This work was supported by UIUC Center for Computational Biotechnology and Genomic Medicine (grant: CSN 1624790)

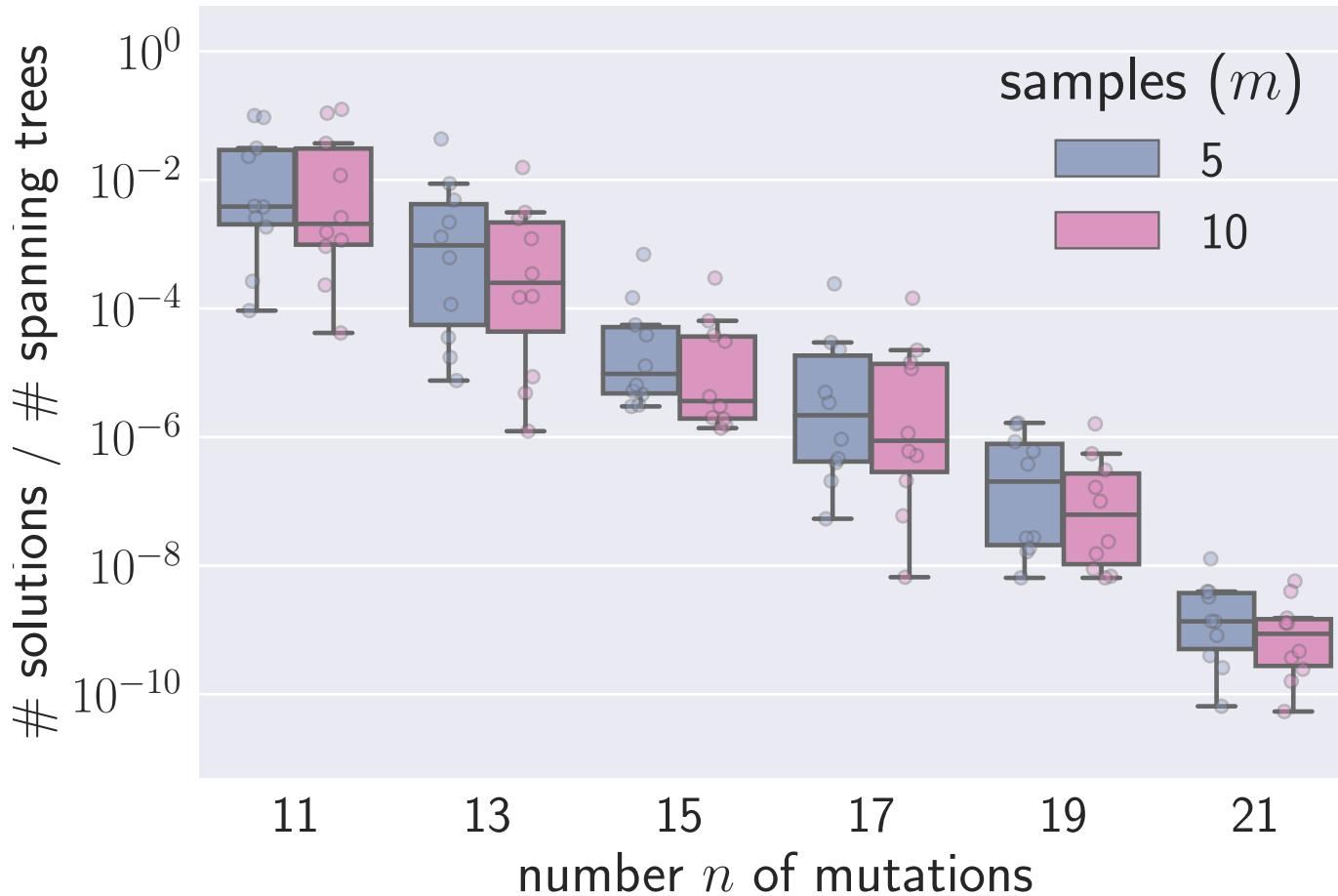
An Upper Bound for Number of Solutions



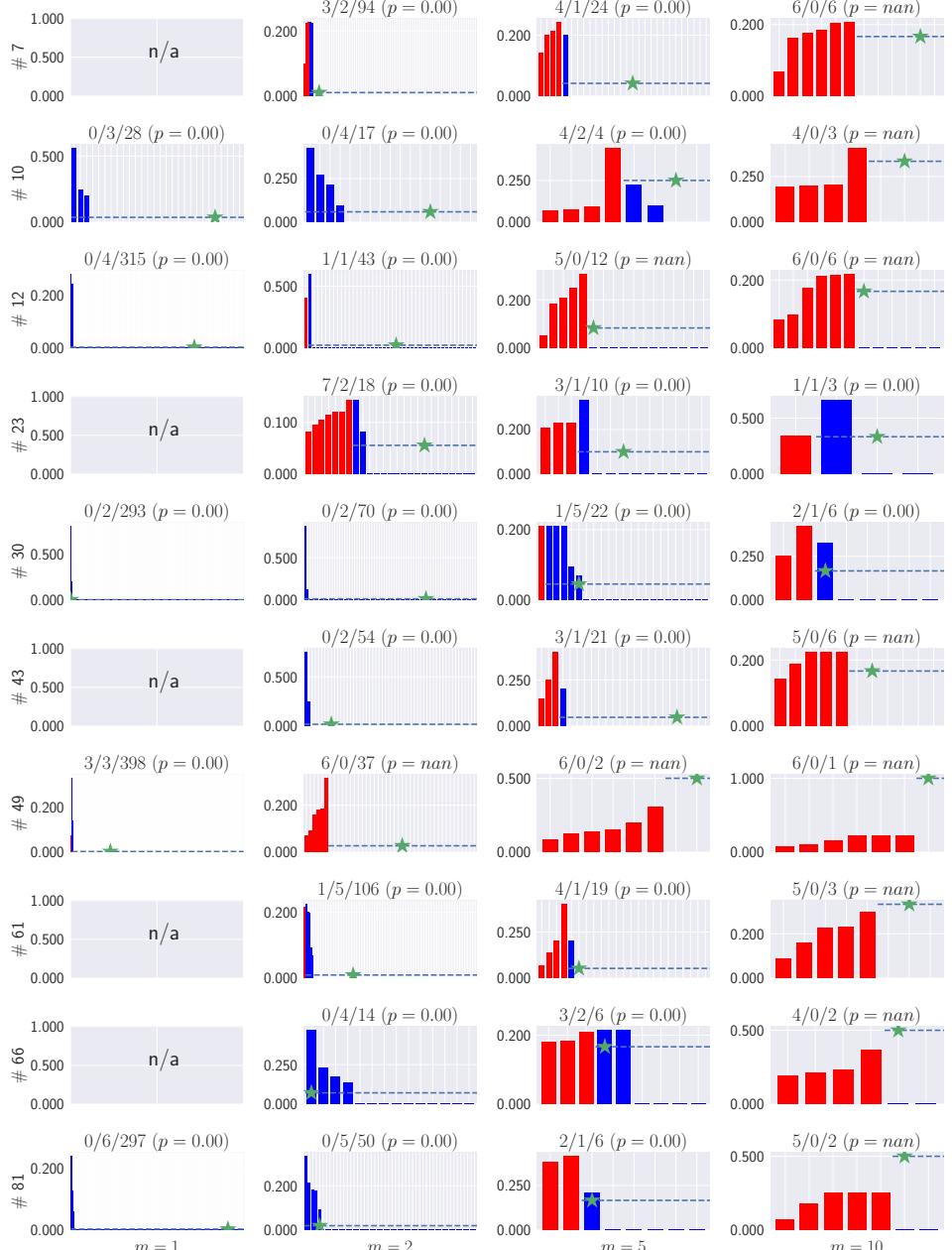
An Upper Bound for Number of Solutions



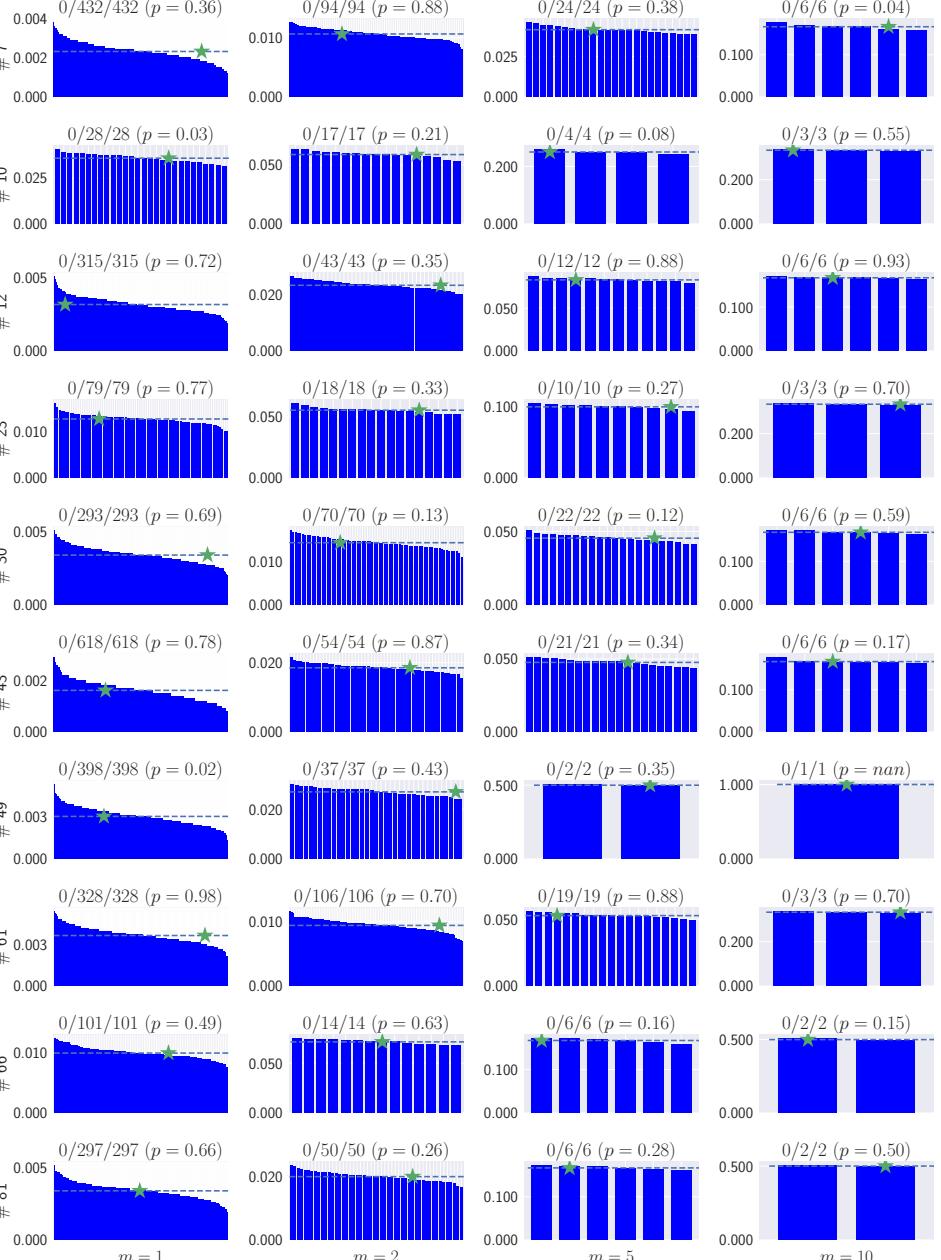
Rejection Sampling Does Not Scale



Canopy



Rejection Sampling



Somatic Mutations Occur at Different Genomic Scales

