

On the Non-uniqueness of Solutions to the Perfect Phylogeny Mixture Problem

Dikshant Pradhan and Mohammed El-Kebir

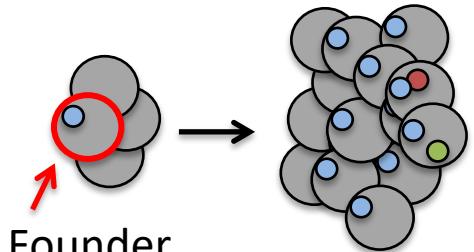
University of Illinois at Urbana Champaign,
Department of Computer Science

RECOMB-CG 2018



Tumorigenesis: Cell Mutation

Clonal Evolution Theory of Cancer
[Nowell, 1976]

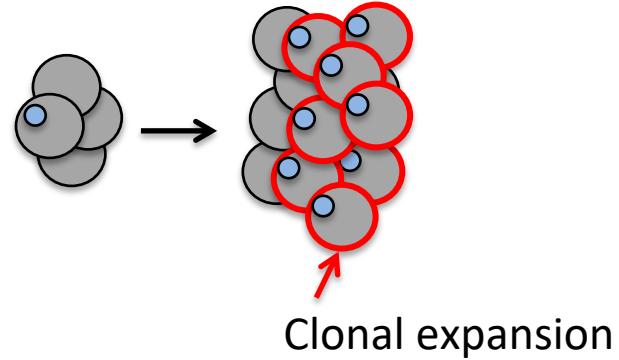


Founder
tumor cell

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

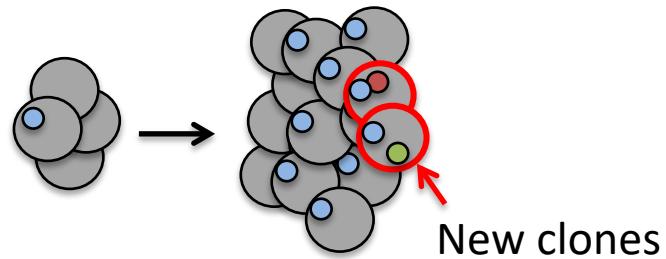
Tumorigenesis: Cell Mutation

Clonal Evolution Theory of Cancer
[Nowell, 1976]



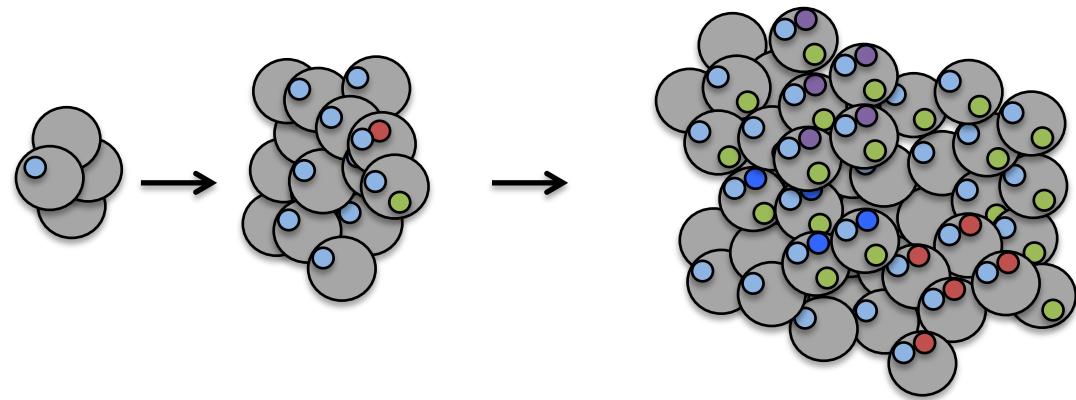
Tumorigenesis: Cell Mutation

Clonal Evolution Theory of Cancer
[Nowell, 1976]



Tumorigenesis: Cell Mutation & Division

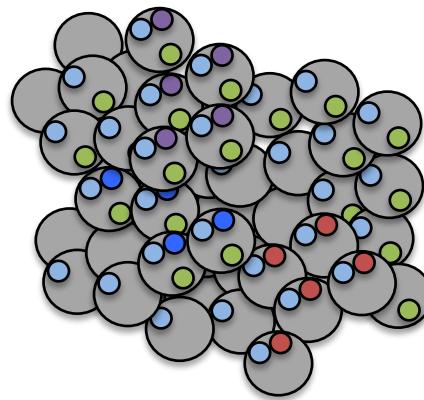
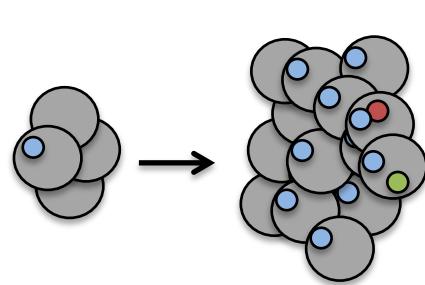
Clonal Evolution Theory of Cancer
[Nowell, 1976]



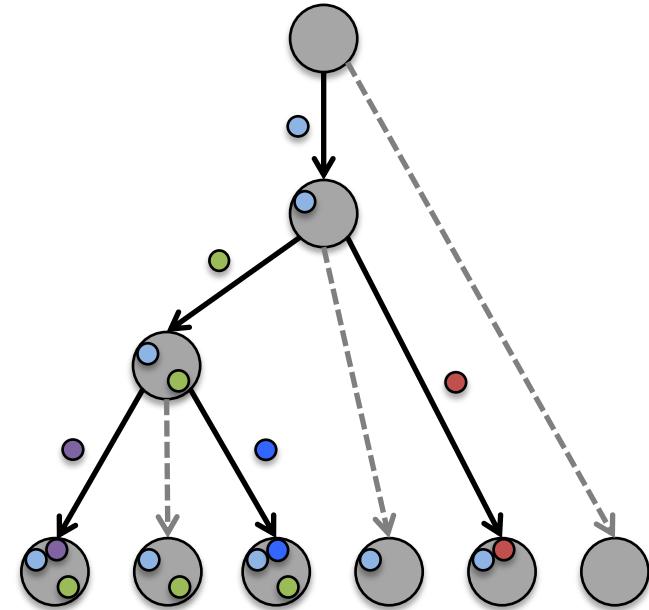
Intra-Tumor
Heterogeneity

Tumorigenesis: Cell Mutation & Division

Clonal Evolution Theory of Cancer
[Nowell, 1976]



Intra-Tumor
Heterogeneity

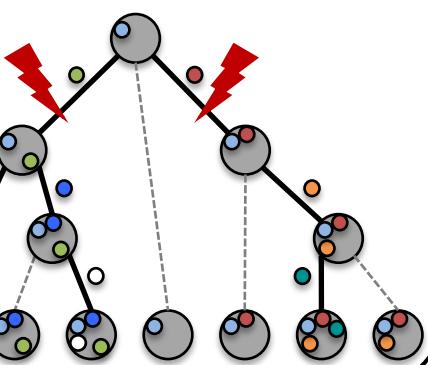


Phylogenetic
Tree T

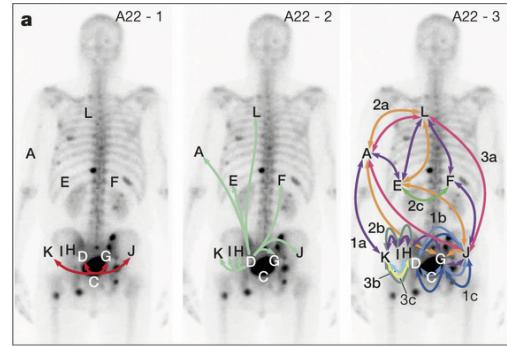
Question: Why are tumor phylogenies important?

Phylogenies are Key to Understanding Cancer

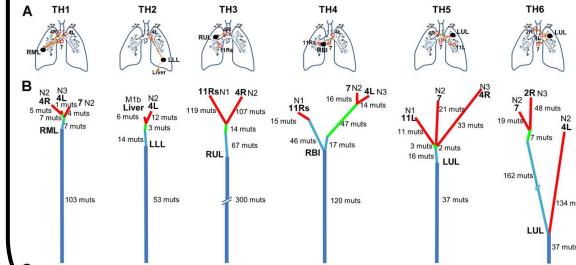
Identify targets for treatment



Understand metastatic development

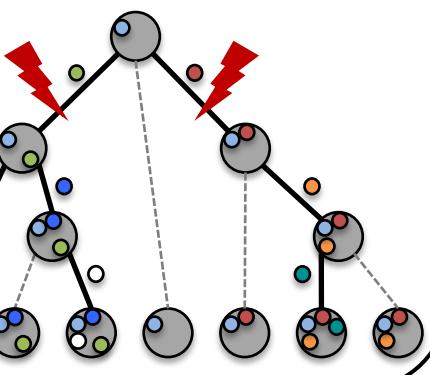


Recognize common patterns of tumor evolution across patients

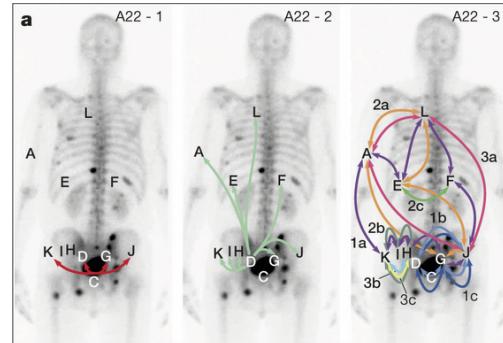


Phylogenies are Key to Understanding Cancer

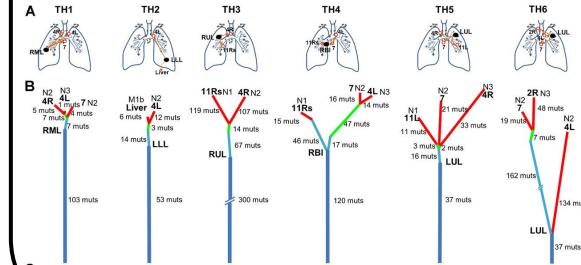
Identify targets for treatment



Understand metastatic development



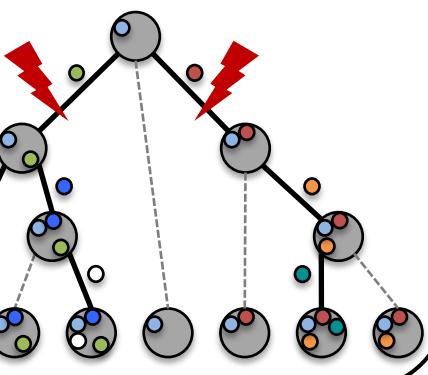
Recognize common patterns of tumor evolution across patients



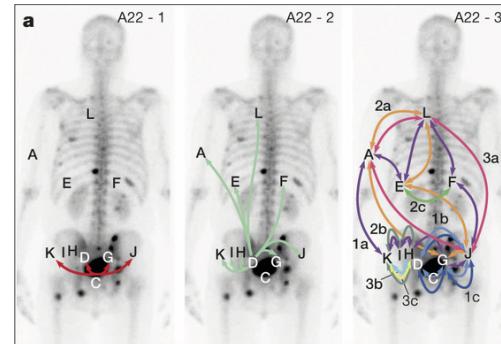
These downstream analyses **critically rely** on accurate tumor phylogeny inference

Phylogenies are Key to Understanding Cancer

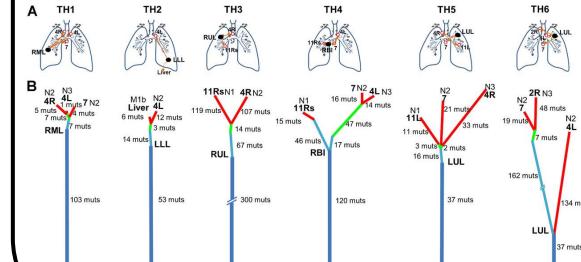
Identify targets for treatment



Understand metastatic development



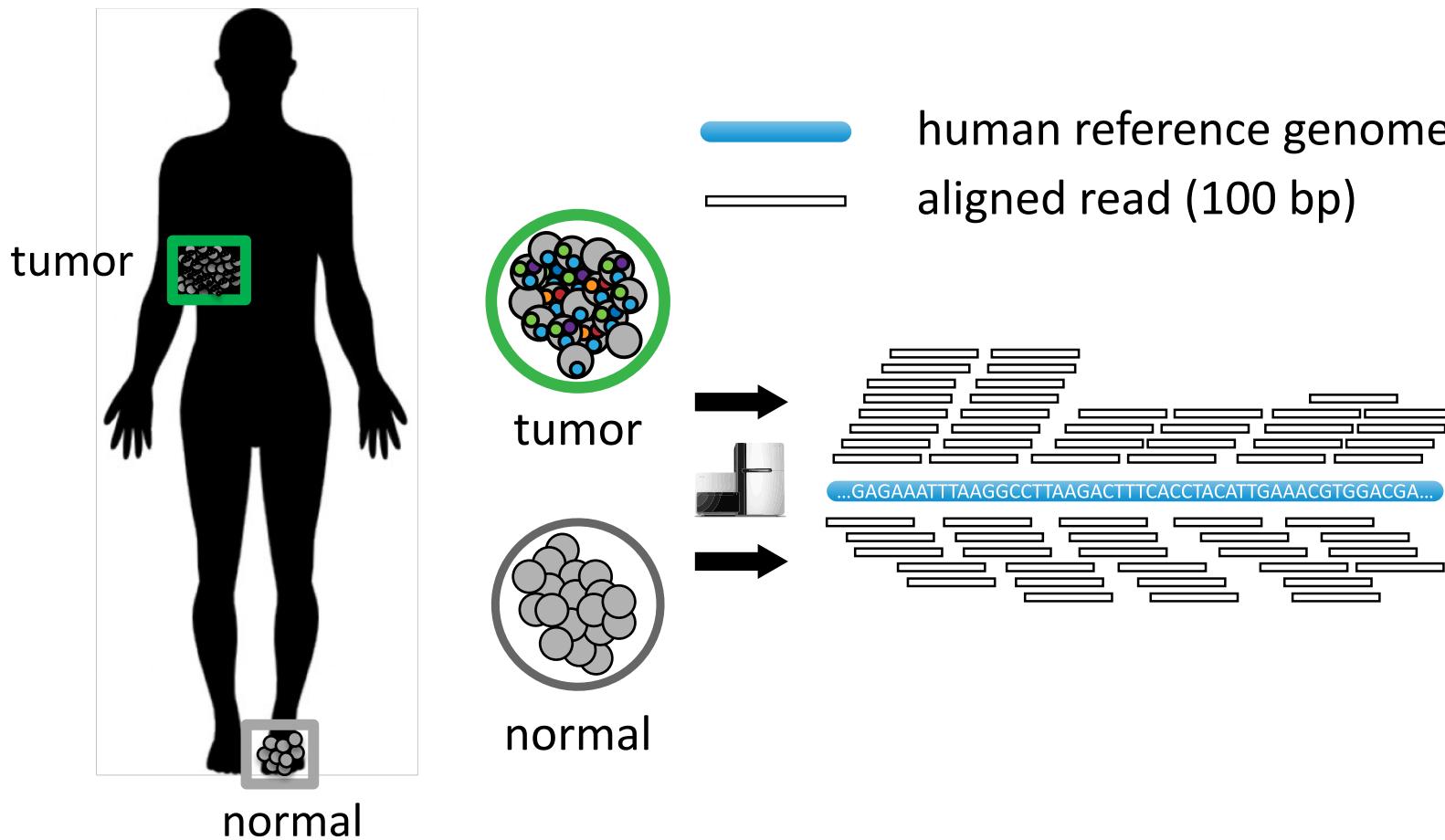
Recognize common patterns of tumor evolution across patients



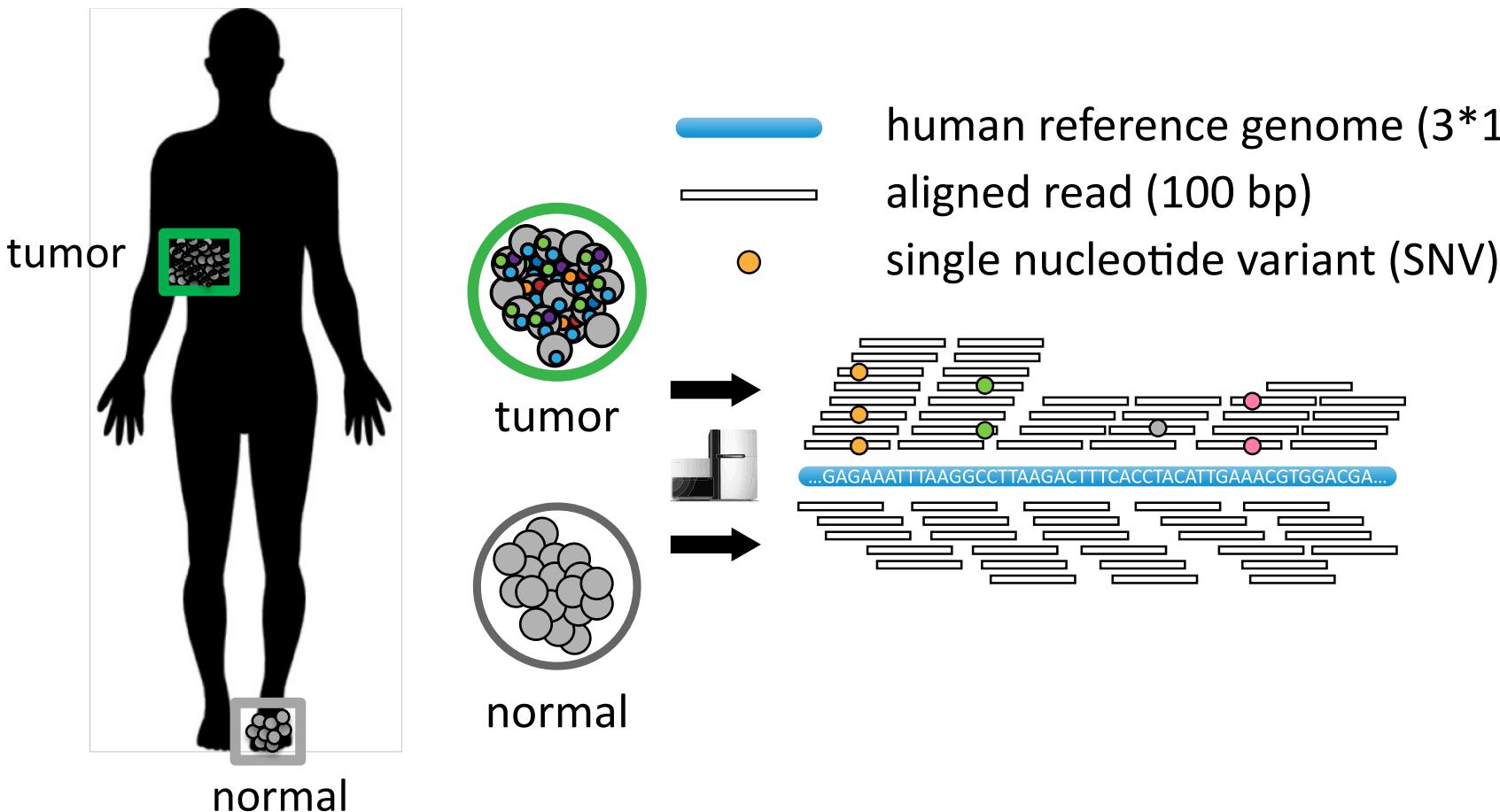
These downstream analyses **critically rely** on accurate tumor phylogeny inference

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

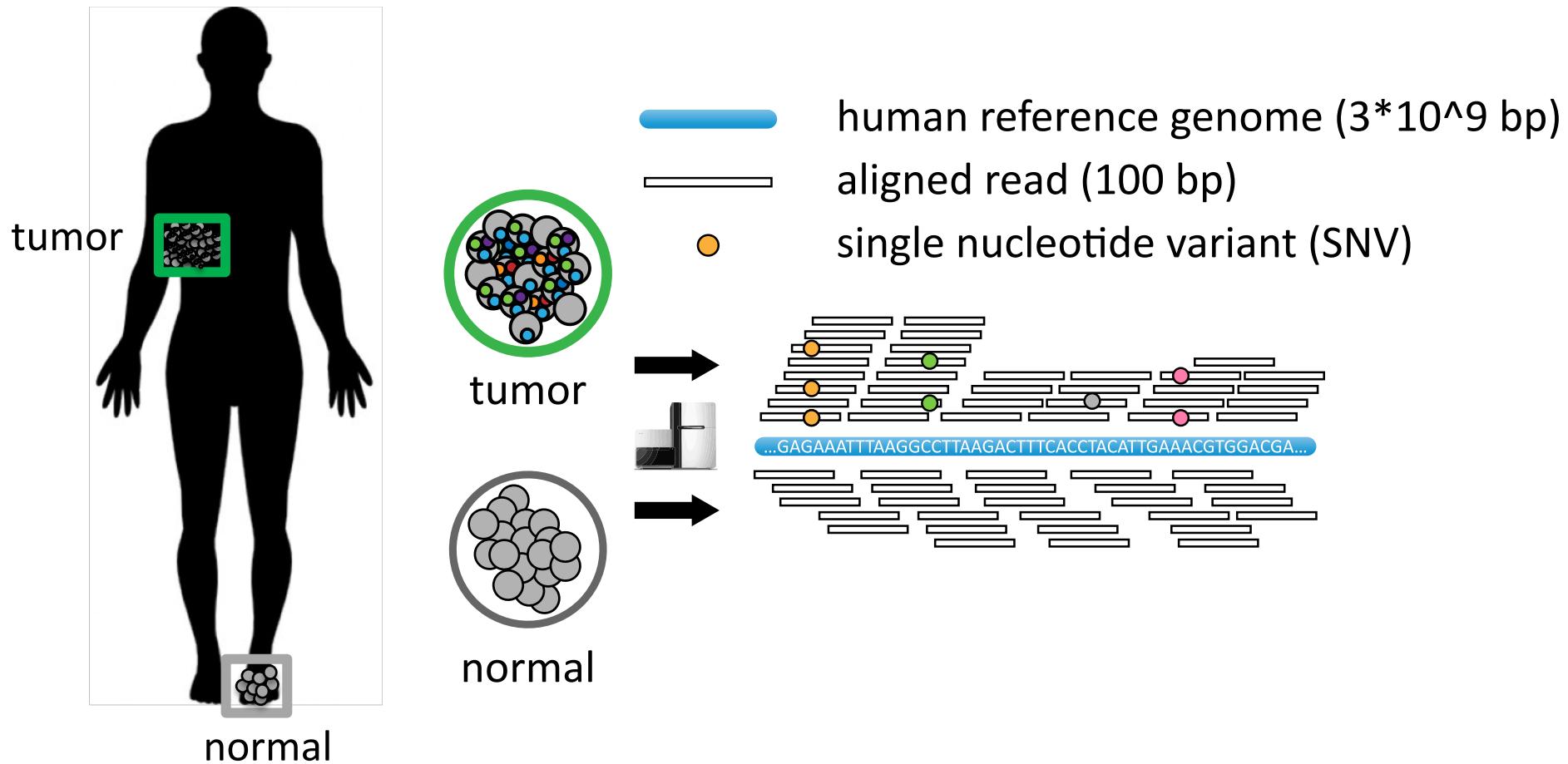
Additional Challenge in Cancer Phylogenetics



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

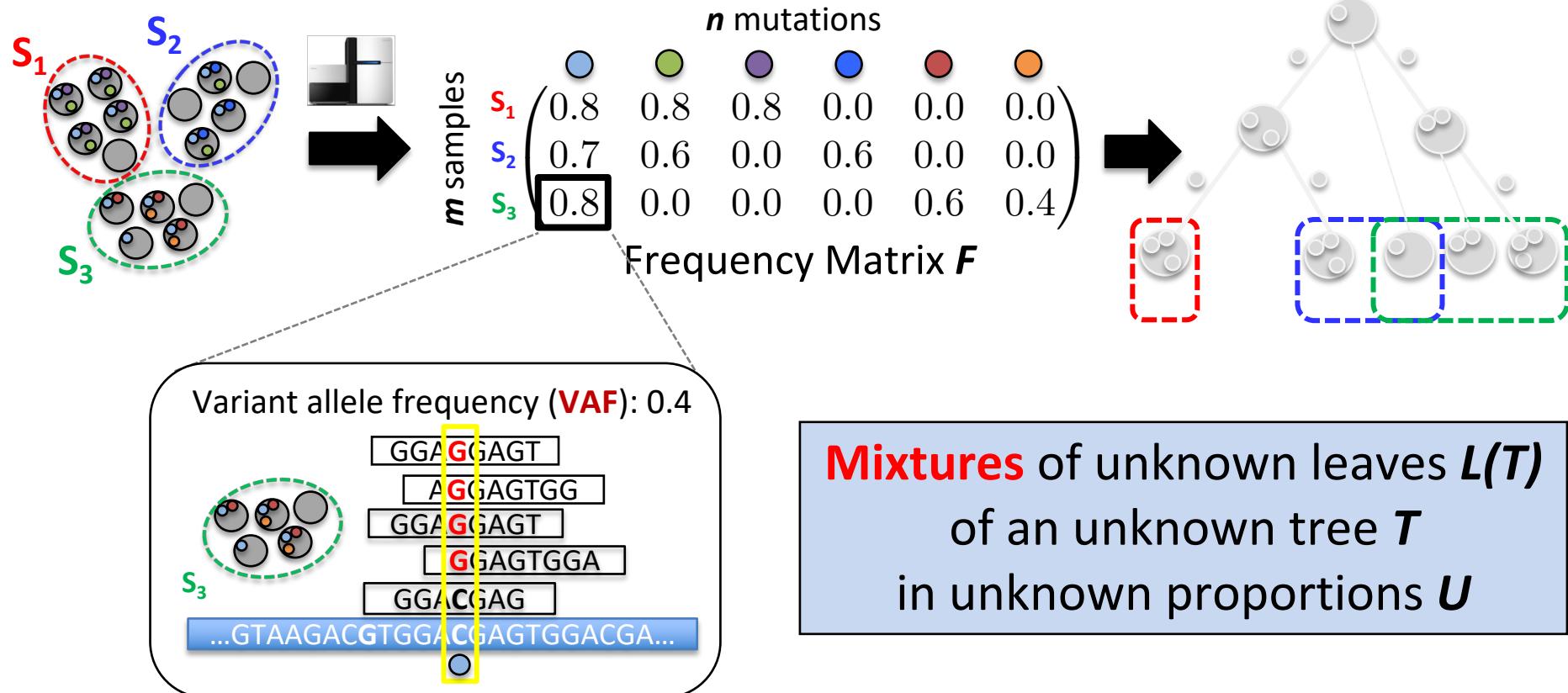
Background and theory:

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

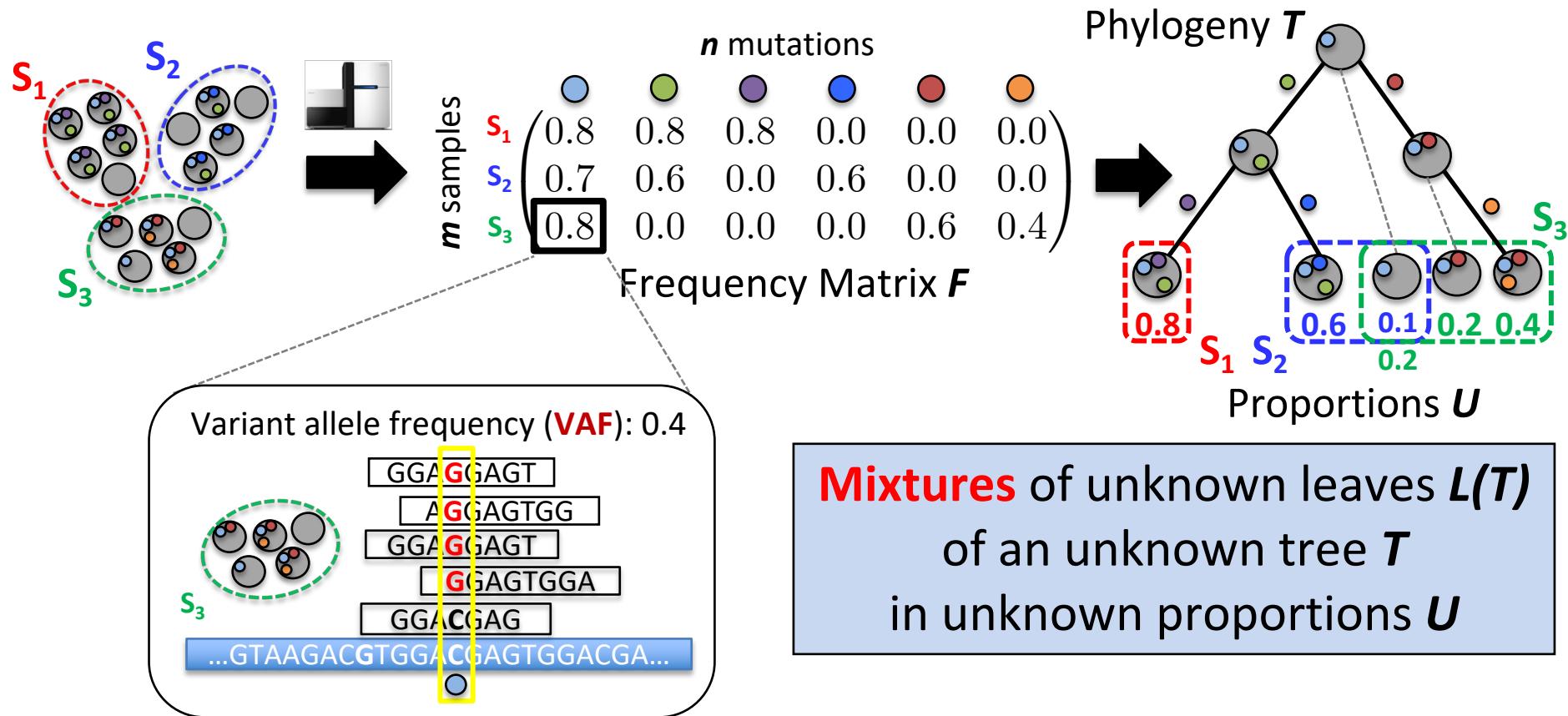
Simulation results:

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

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

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

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 $\mathbf{F} = \mathbf{U} \mathbf{B}$

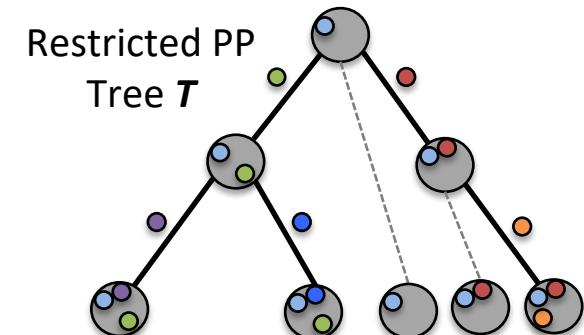
Previous Work

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:

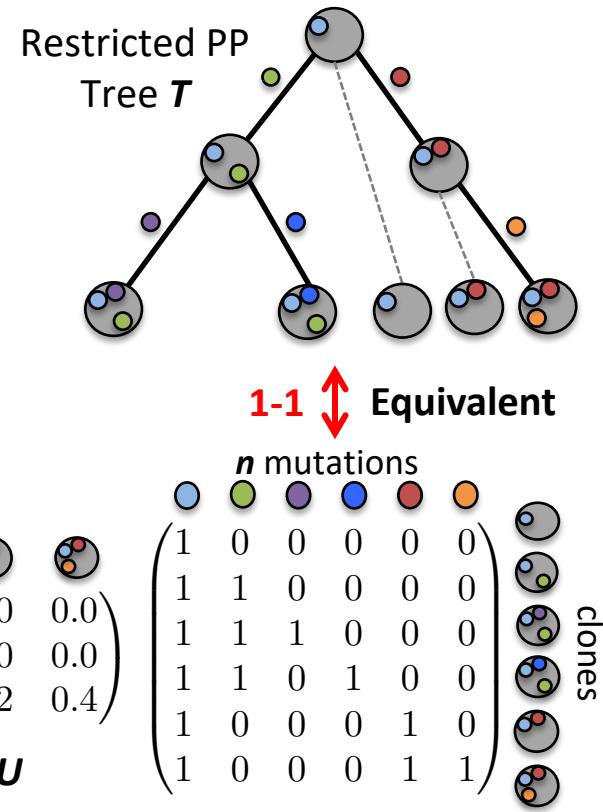
$$u_{pj} \geq 0 \text{ and } \sum_j u_{pj} \leq 1$$

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



Rows of U are proportions:

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

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

Combinatorial Characterization

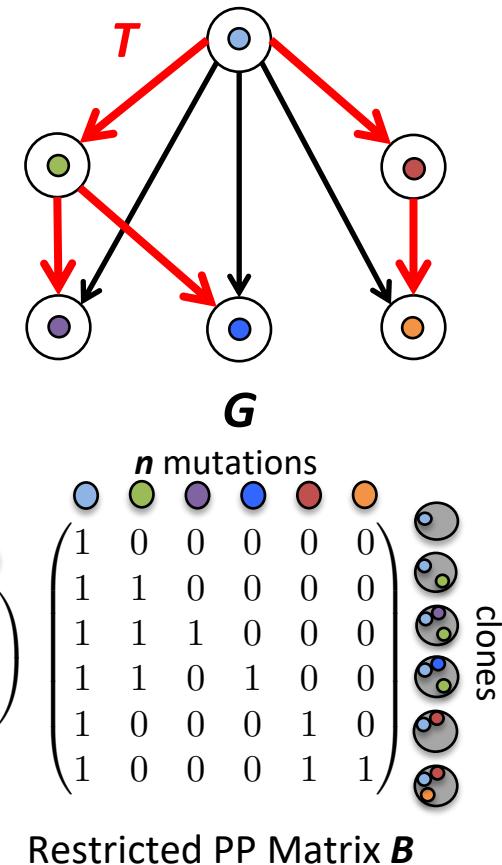
- Frequency $f_{p,i}$ is mass of subtree rooted at node that introduced i
- Usage $u_{p,i}$ is mass of node that introduced i

$$\begin{matrix} m \text{ samples} \\ \text{---} \\ S_1 & S_2 & S_3 \end{matrix} \left(\begin{matrix} n \text{ mutations} \\ 0.8 & 0.8 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.7 & 0.6 & 0.0 & 0.6 & 0.0 & 0.0 \\ 0.8 & 0.0 & 0.0 & 0.0 & 0.6 & 0.4 \end{matrix} \right) = \begin{matrix} m \text{ samples} \\ \text{---} \\ S_1 & S_2 & S_3 \end{matrix} \left(\begin{matrix} n \text{ mutations} \\ 0.0 & 0.0 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.1 & 0.0 & 0.0 & 0.6 & 0.0 & 0.0 \\ 0.2 & 0.0 & 0.0 & 0.0 & 0.2 & 0.4 \end{matrix} \right)$$

Frequency Matrix \mathbf{F}

$$= \begin{matrix} m \text{ samples} \\ \text{---} \\ S_1 & S_2 & S_3 \end{matrix} \left(\begin{matrix} n \text{ mutations} \\ 0.0 & 0.0 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.1 & 0.0 & 0.0 & 0.6 & 0.0 & 0.0 \\ 0.2 & 0.0 & 0.0 & 0.0 & 0.2 & 0.4 \end{matrix} \right)$$

Mixture Matrix \mathbf{U}



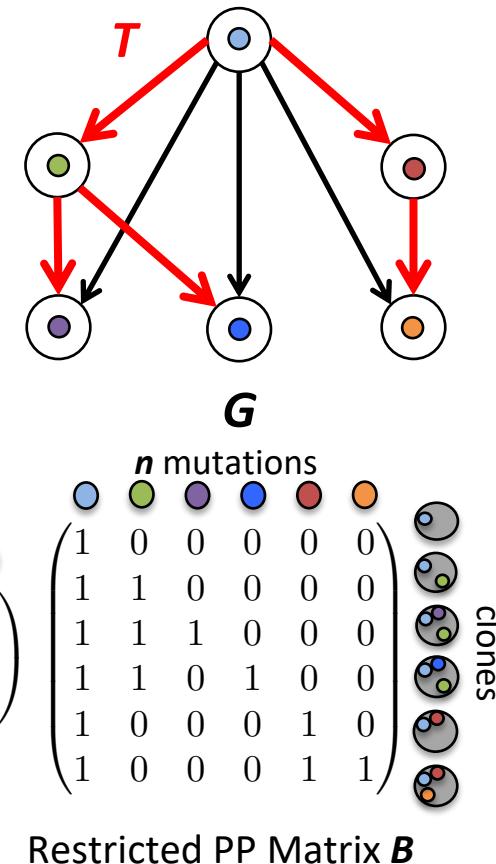
Theorem 1:

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

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



Theorem 1:

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

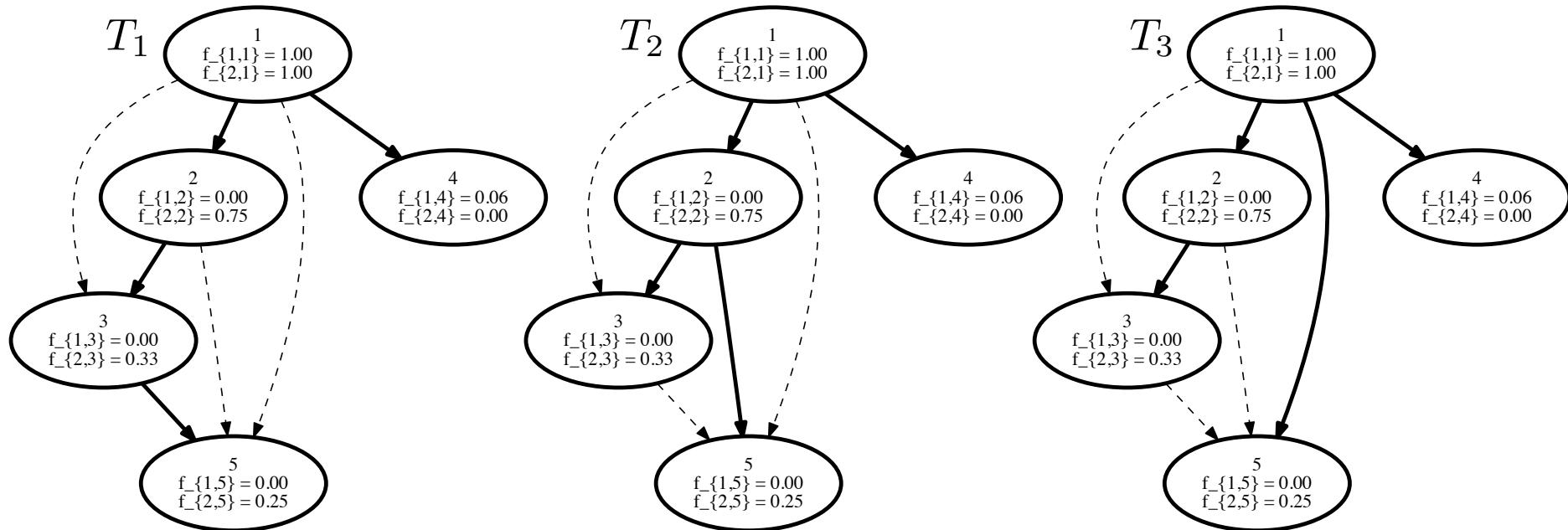
Theorem 2:

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

Perfect Phylogeny Mixture: [El-Kebir*, Oesper* et al., 2015]

Given F , find U and B such that $F = UB$

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

Background and theory:

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

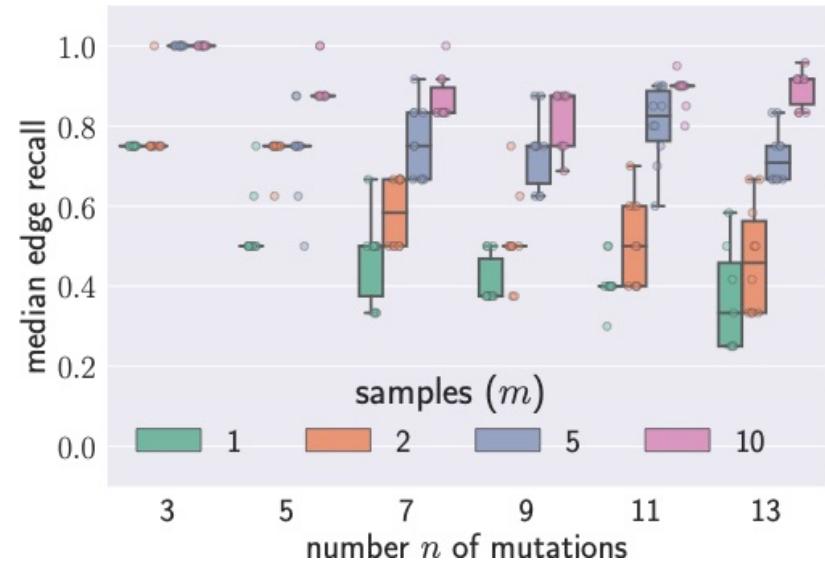
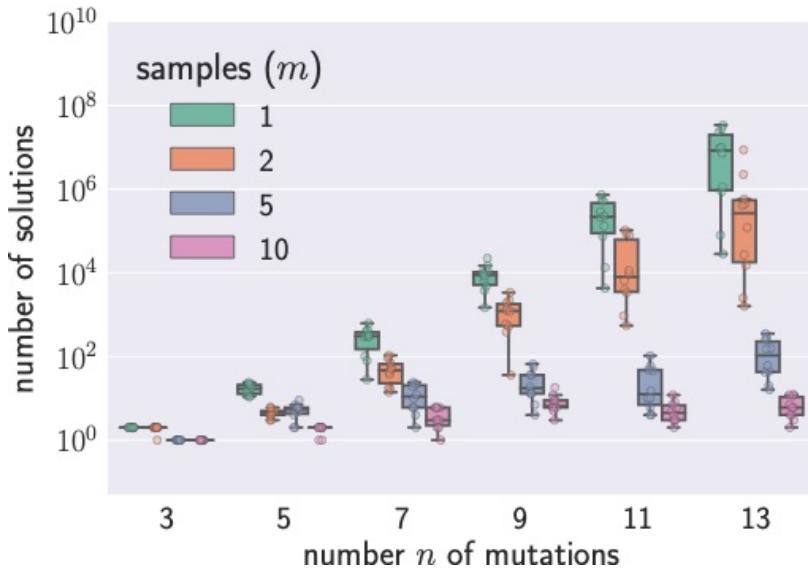
Simulation results:

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

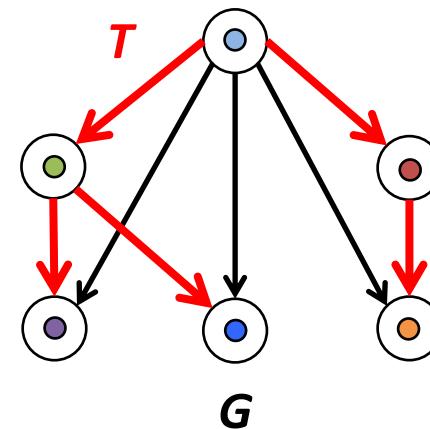
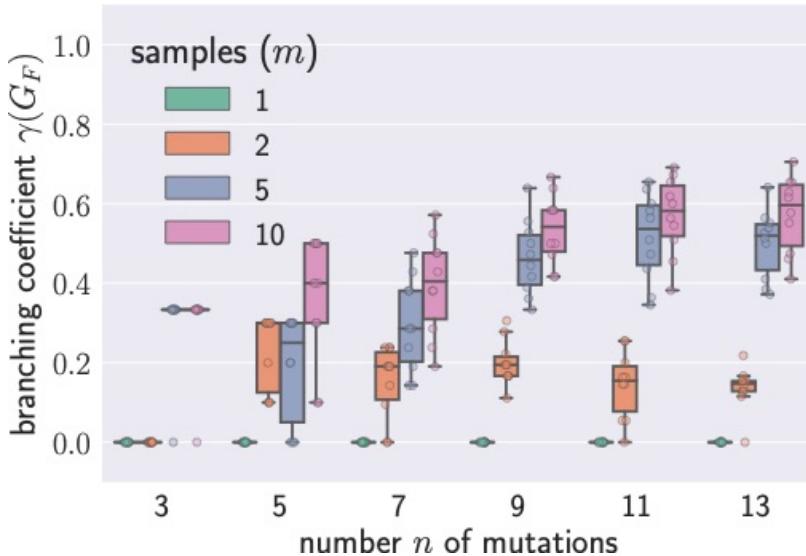
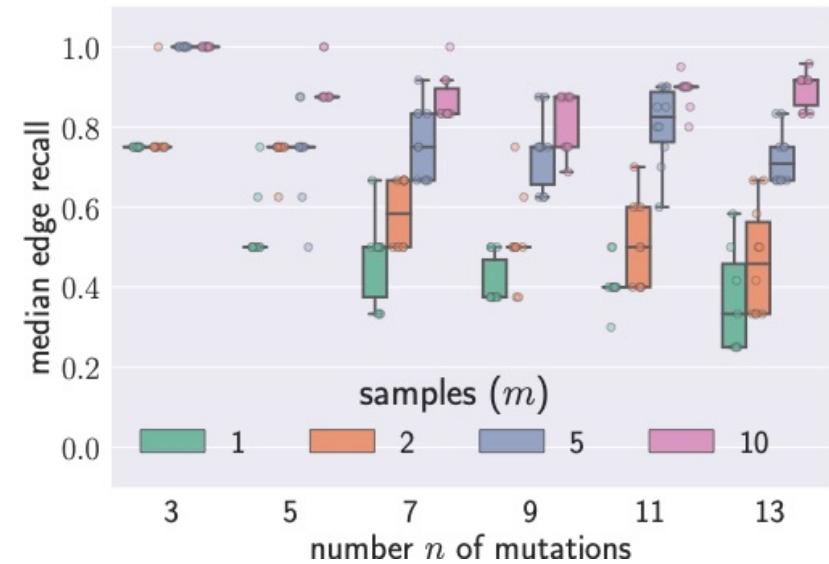
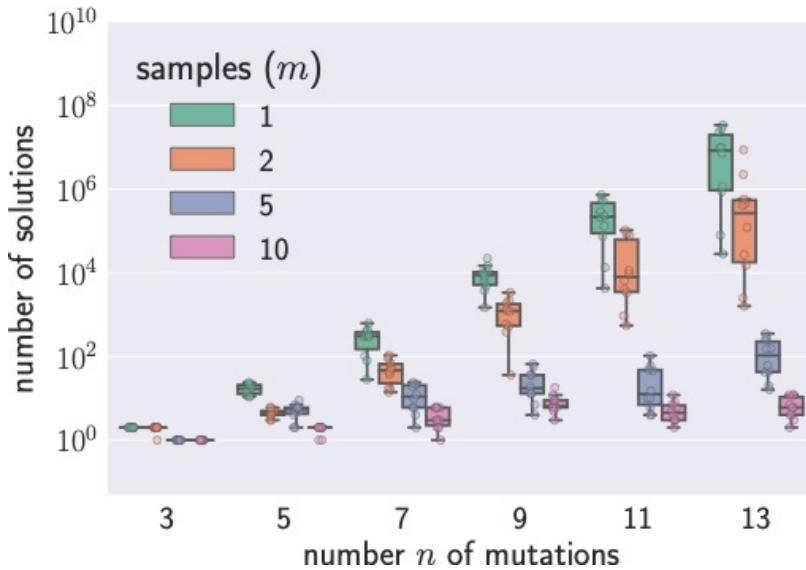


Dikshant Pradhan

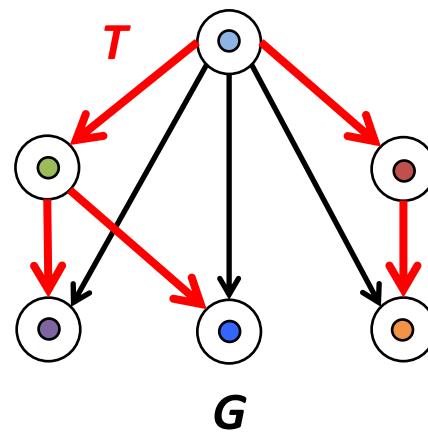
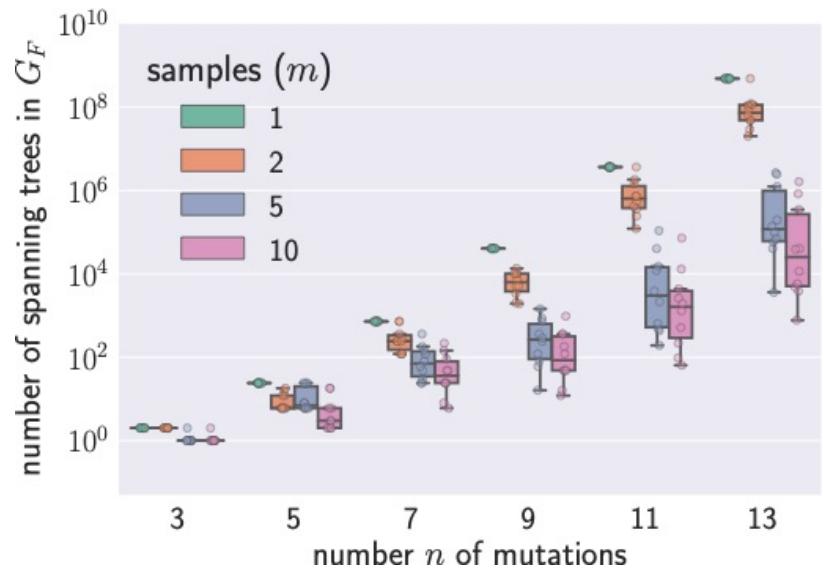
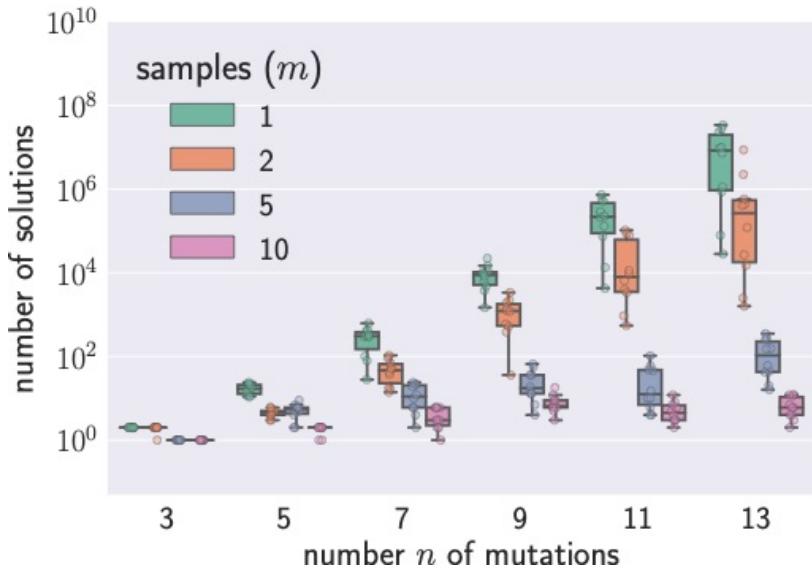
What Contributors to Non-uniqueness?



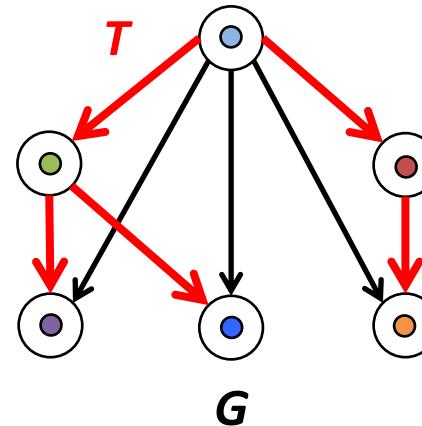
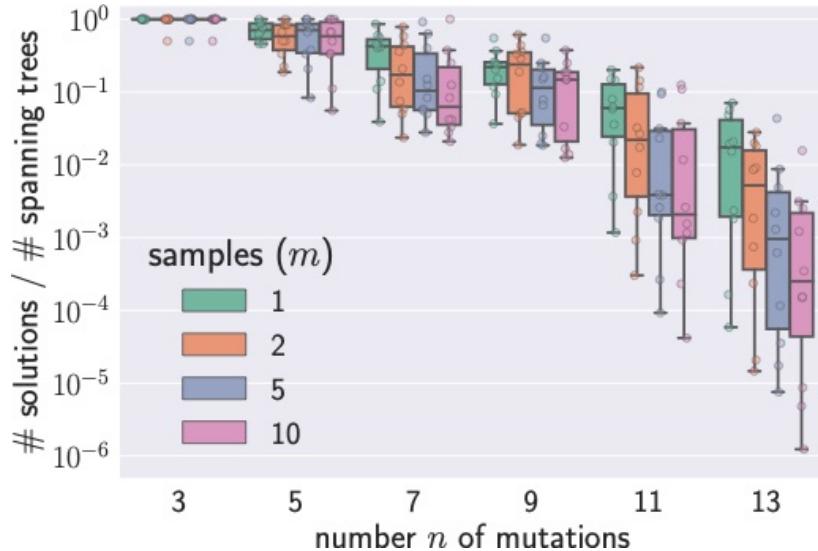
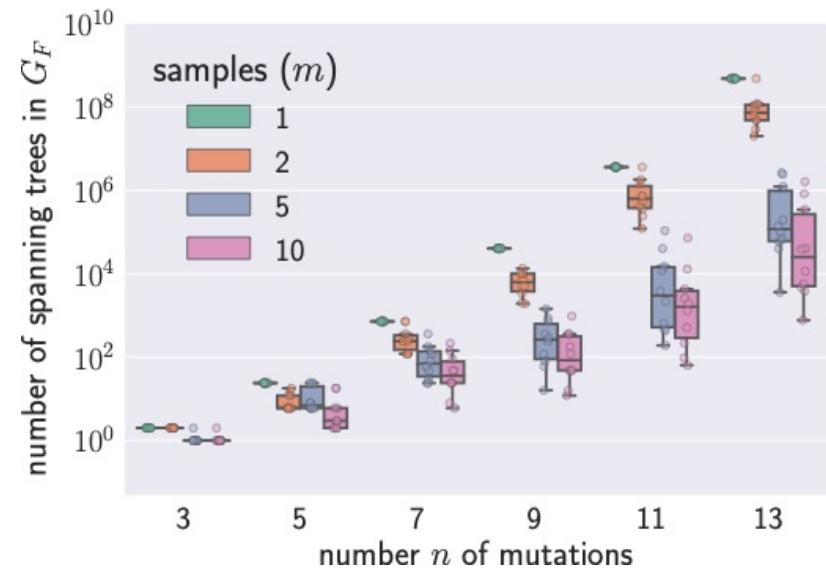
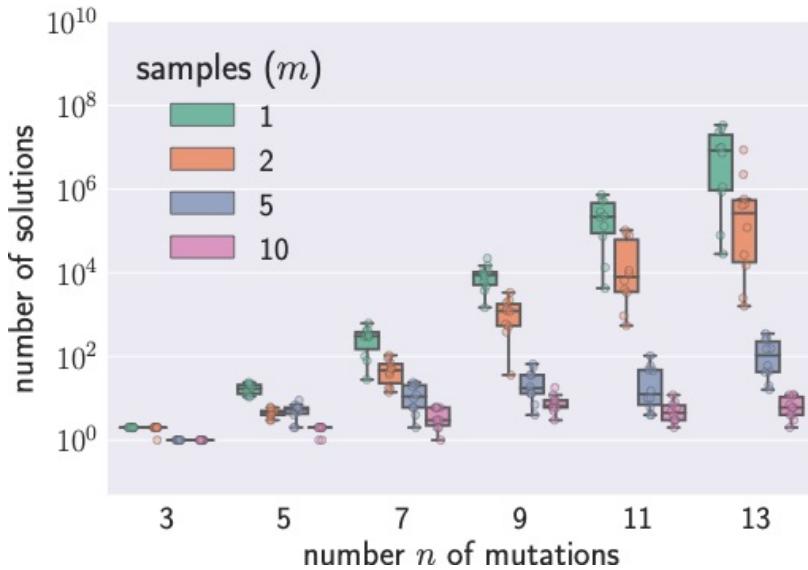
What Contributors to Non-uniqueness?



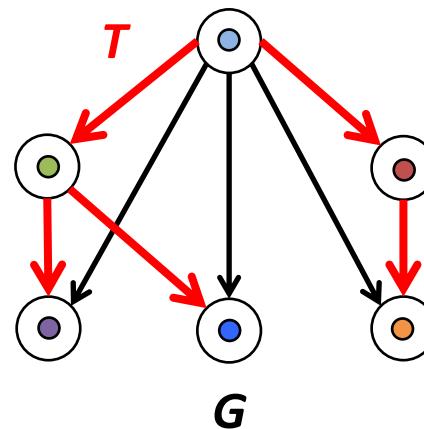
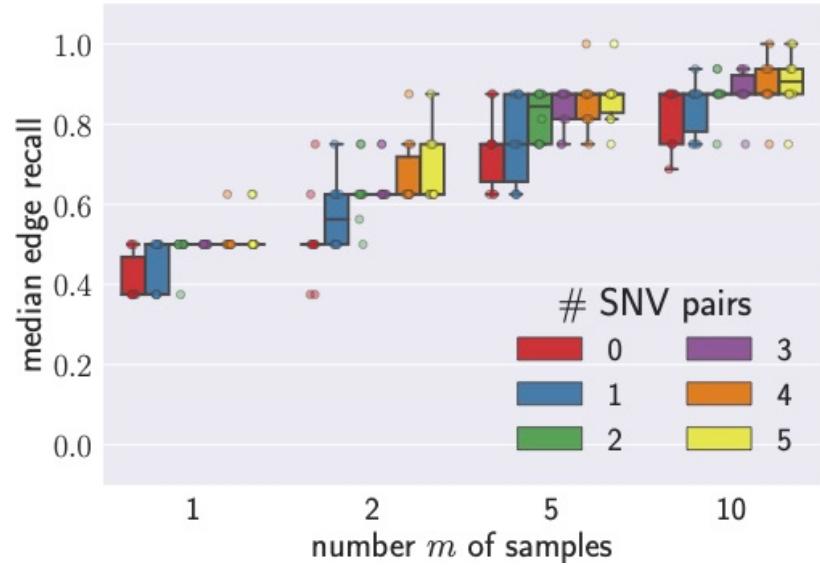
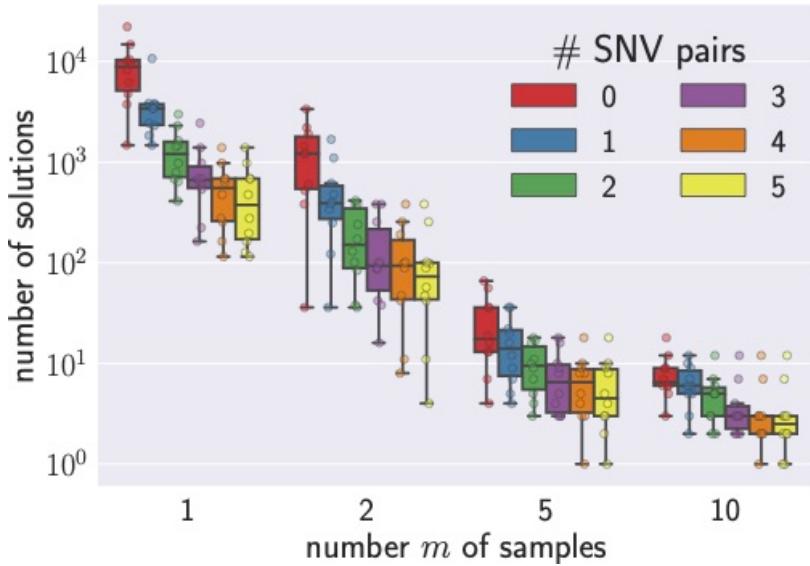
An Upper Bound for Number of Solutions



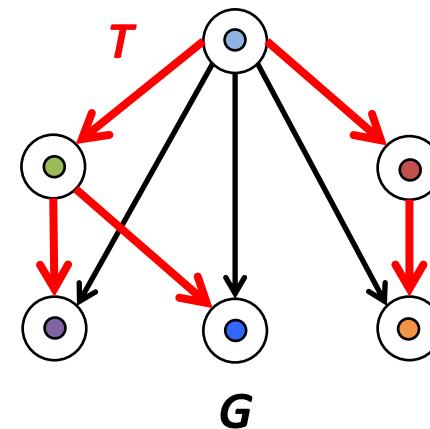
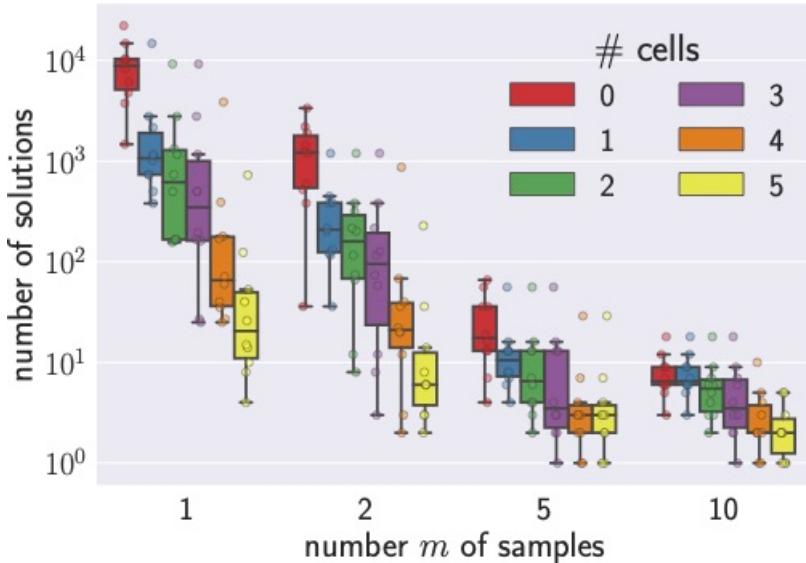
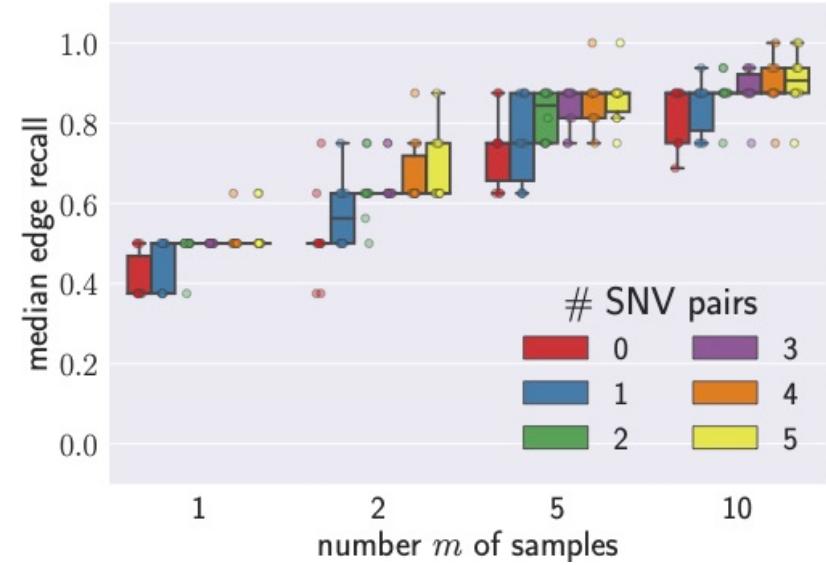
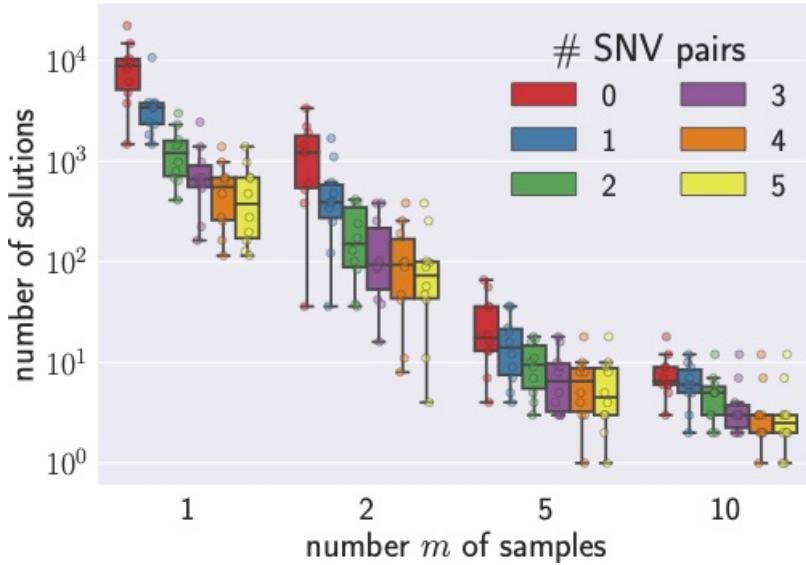
An Upper Bound for Number of Solutions



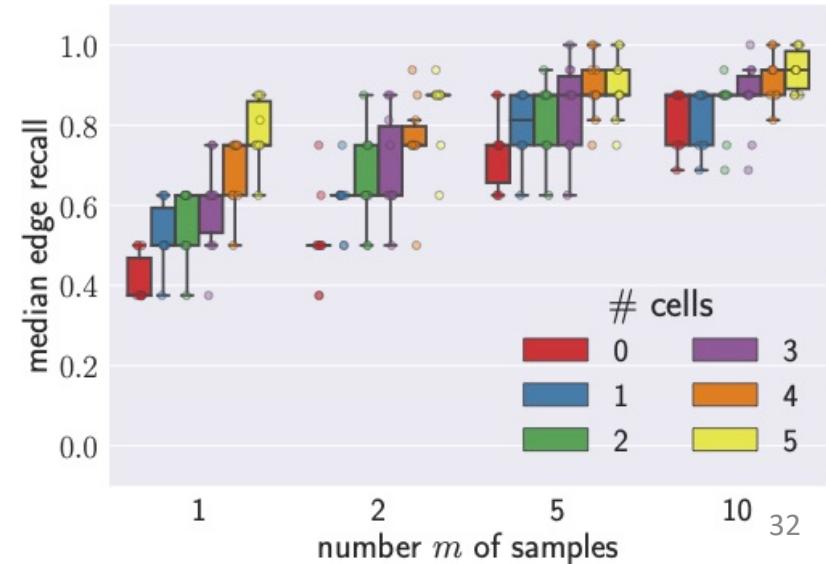
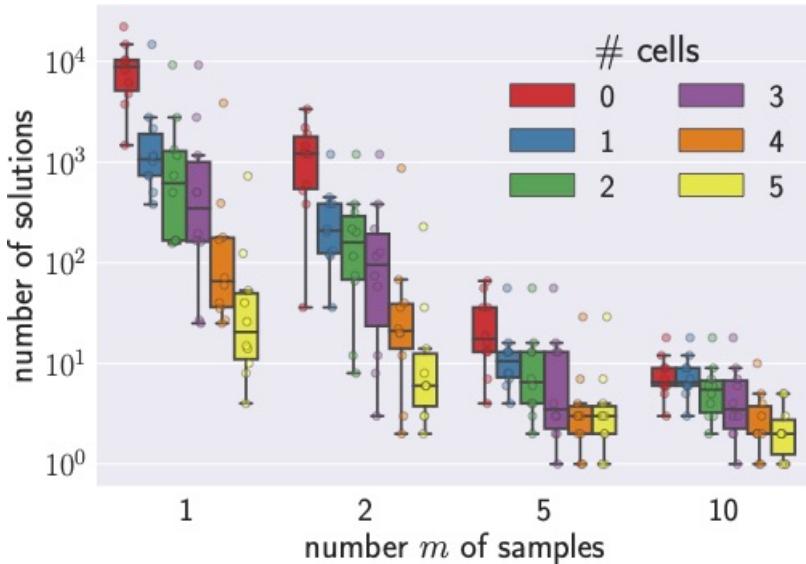
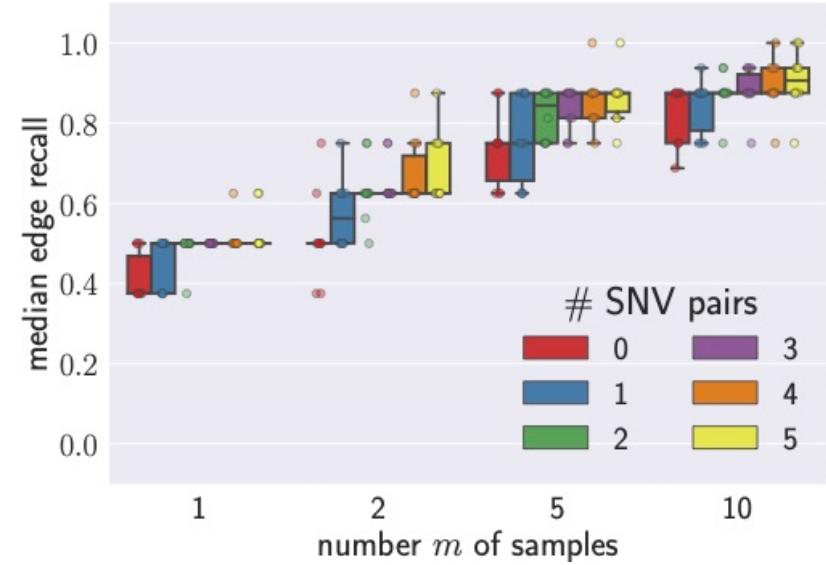
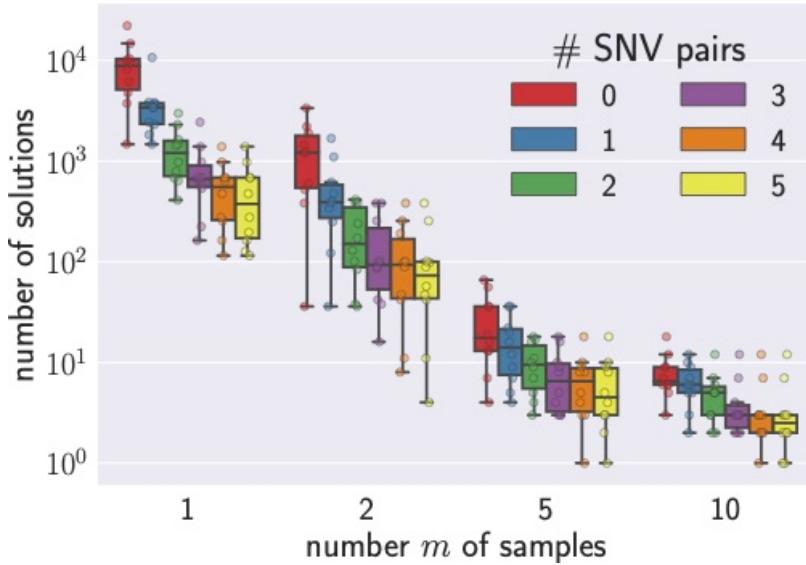
How to Reduce Non-Uniqueness?



How to Reduce Non-Uniqueness?



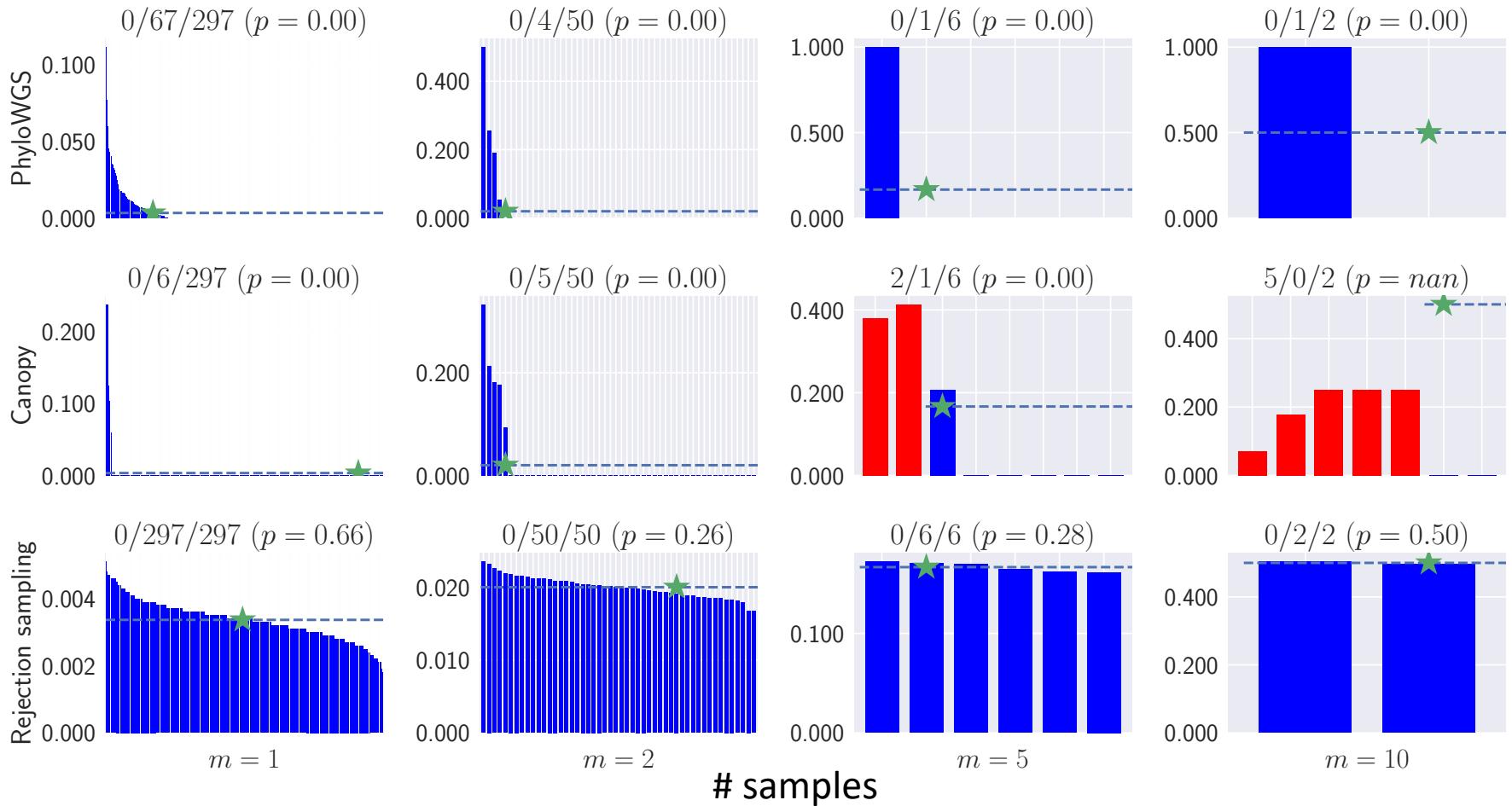
How to Reduce Non-Uniqueness?



How Does Non-uniqueness affect Methods?

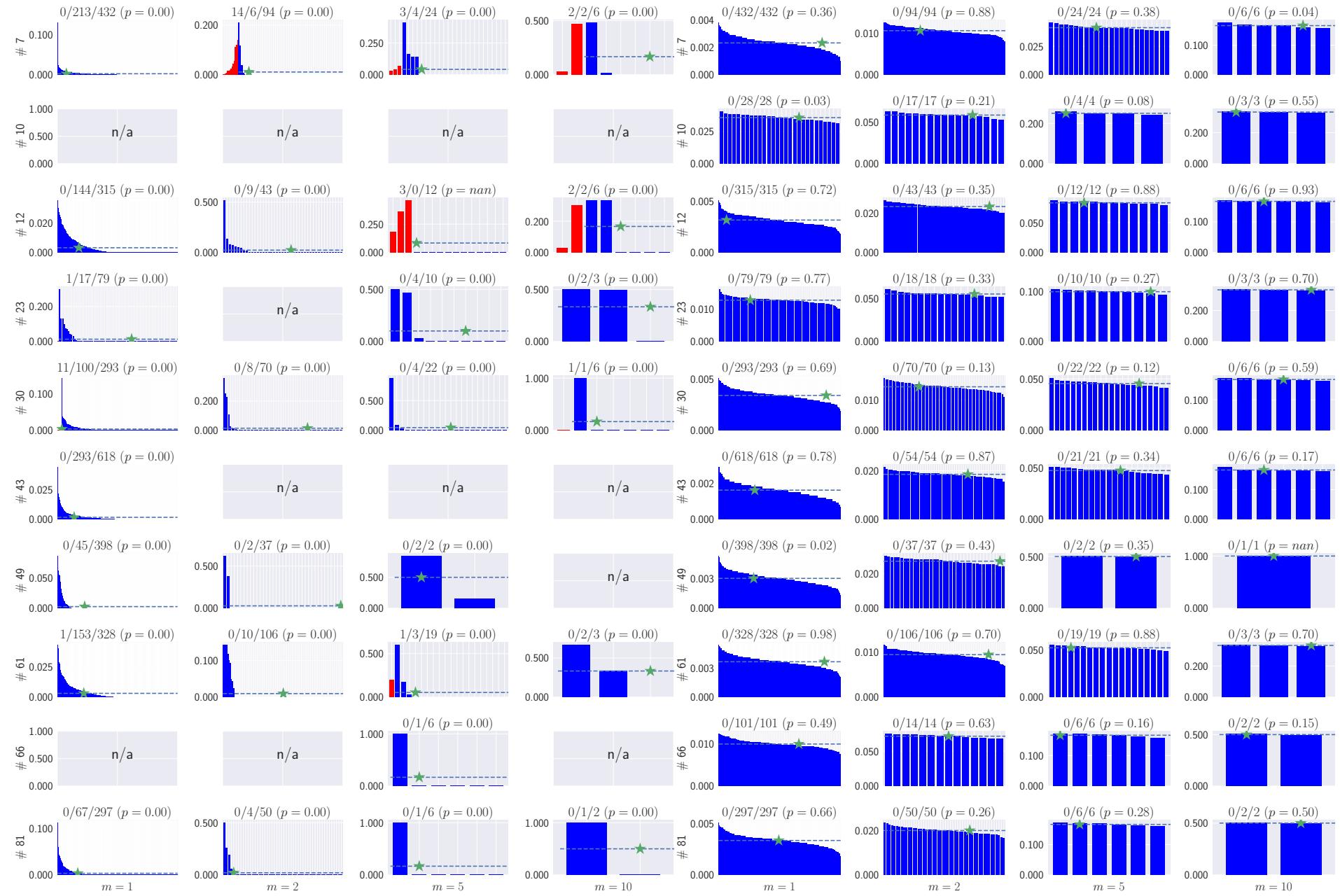
Two current MCMC methods using default parameters:

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

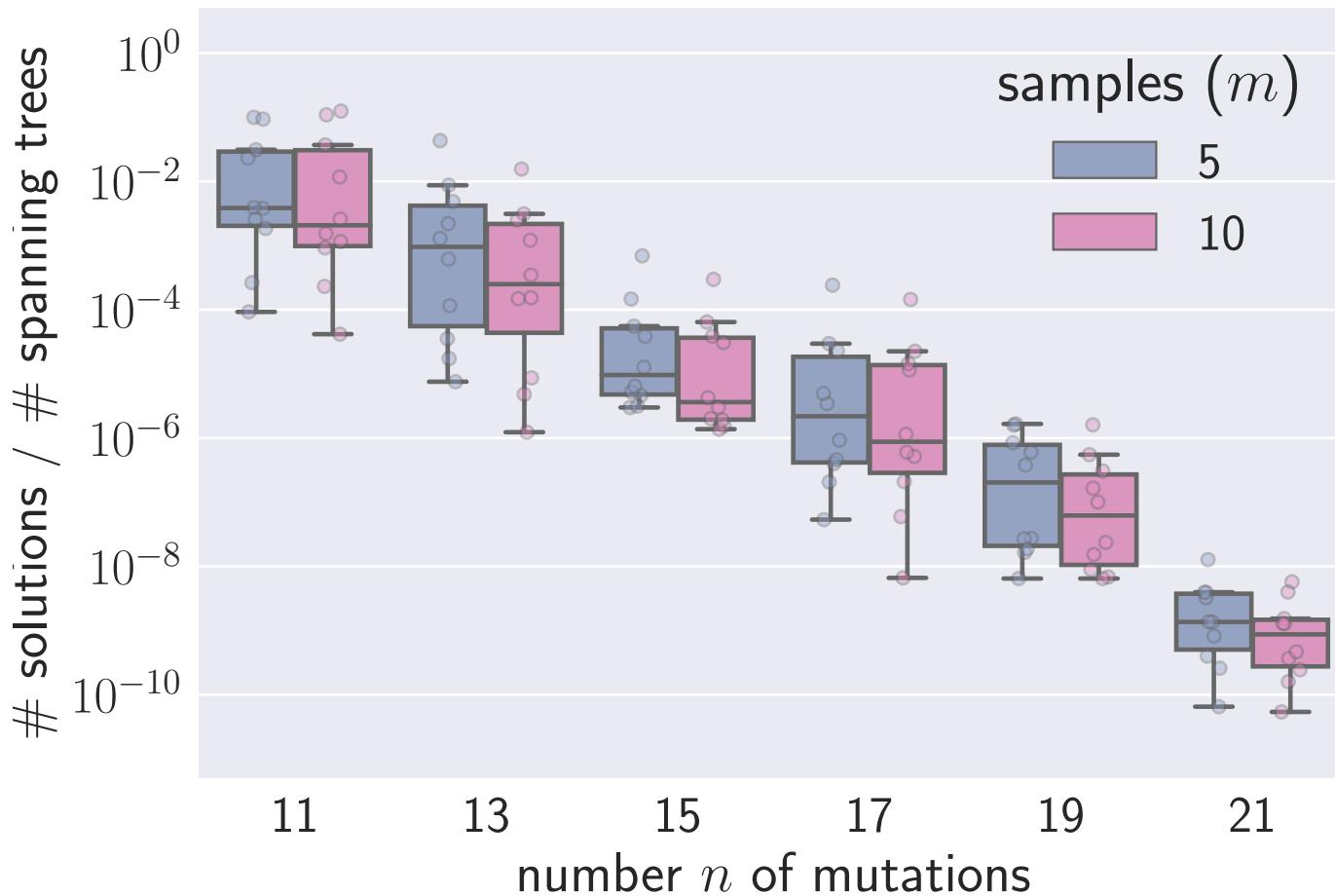


PhyloWGS

Rejection Sampling

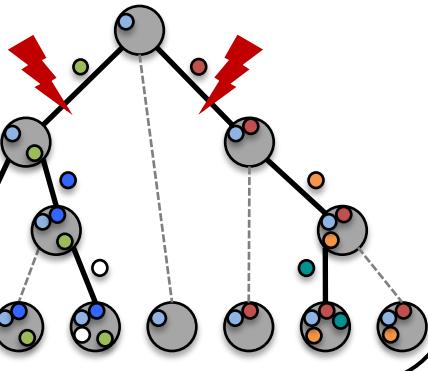


Rejection Sampling Does Not Scale

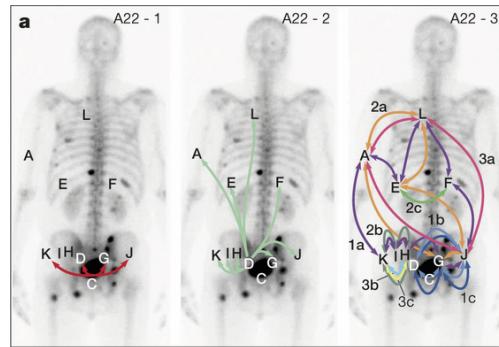


Challenge

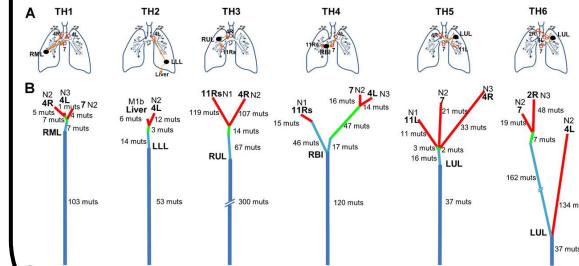
Identify targets for treatment



Understand metastatic development



Recognize common patterns of tumor evolution across patients



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

Key challenge:

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

Conclusion

Background and theory:

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

Simulation results:

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

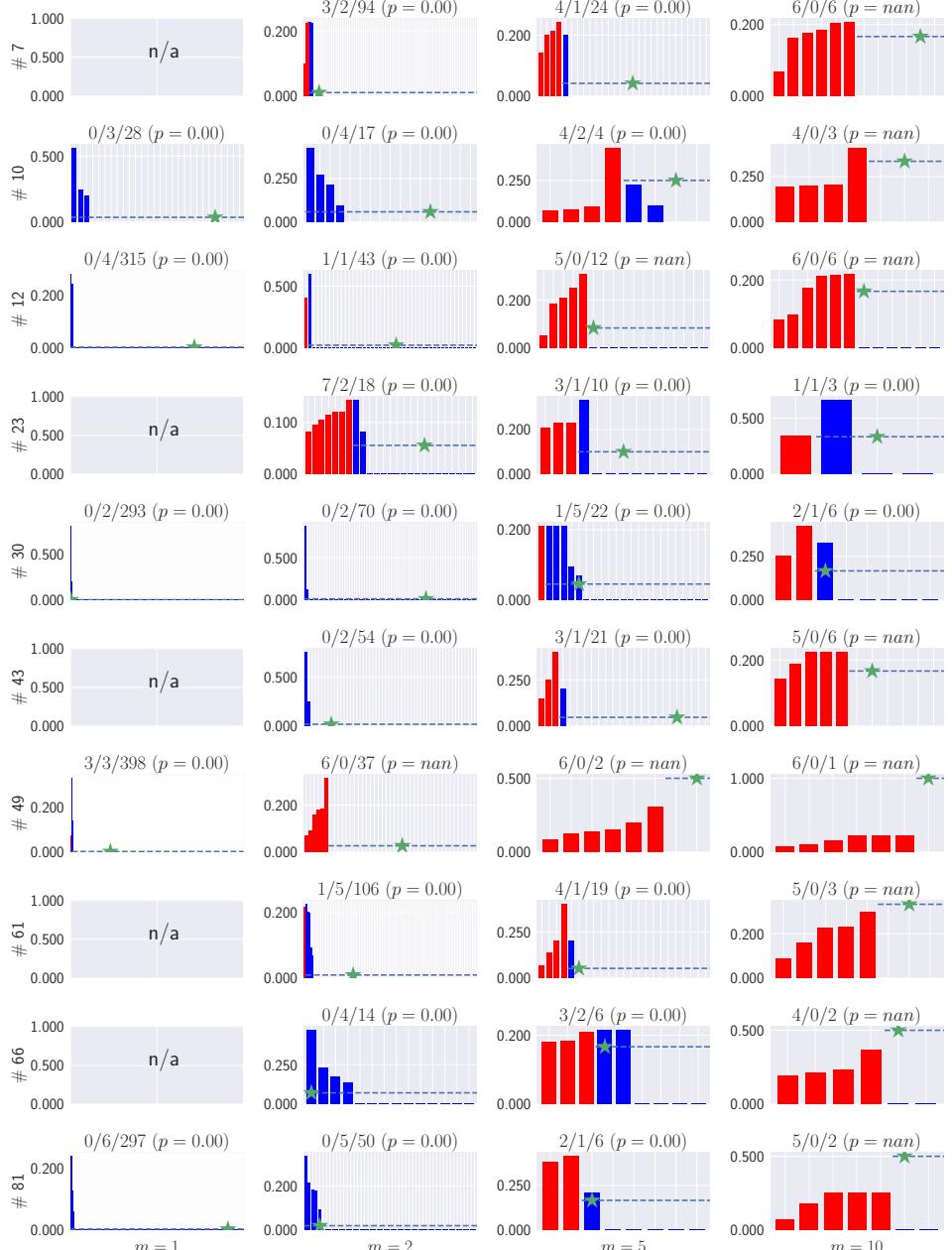
Future Work

- Better sampling and counting algorithms
- Non-uniqueness in an infinite setting: characterize statistical consistency
- Constrained tumor phylogeny inference
 - Metastasis [El-Kebir, Satas & Raphael, Nature Genetics, 2018]
 - SCS + bulk, long-read sequencing
 - Other constraints...

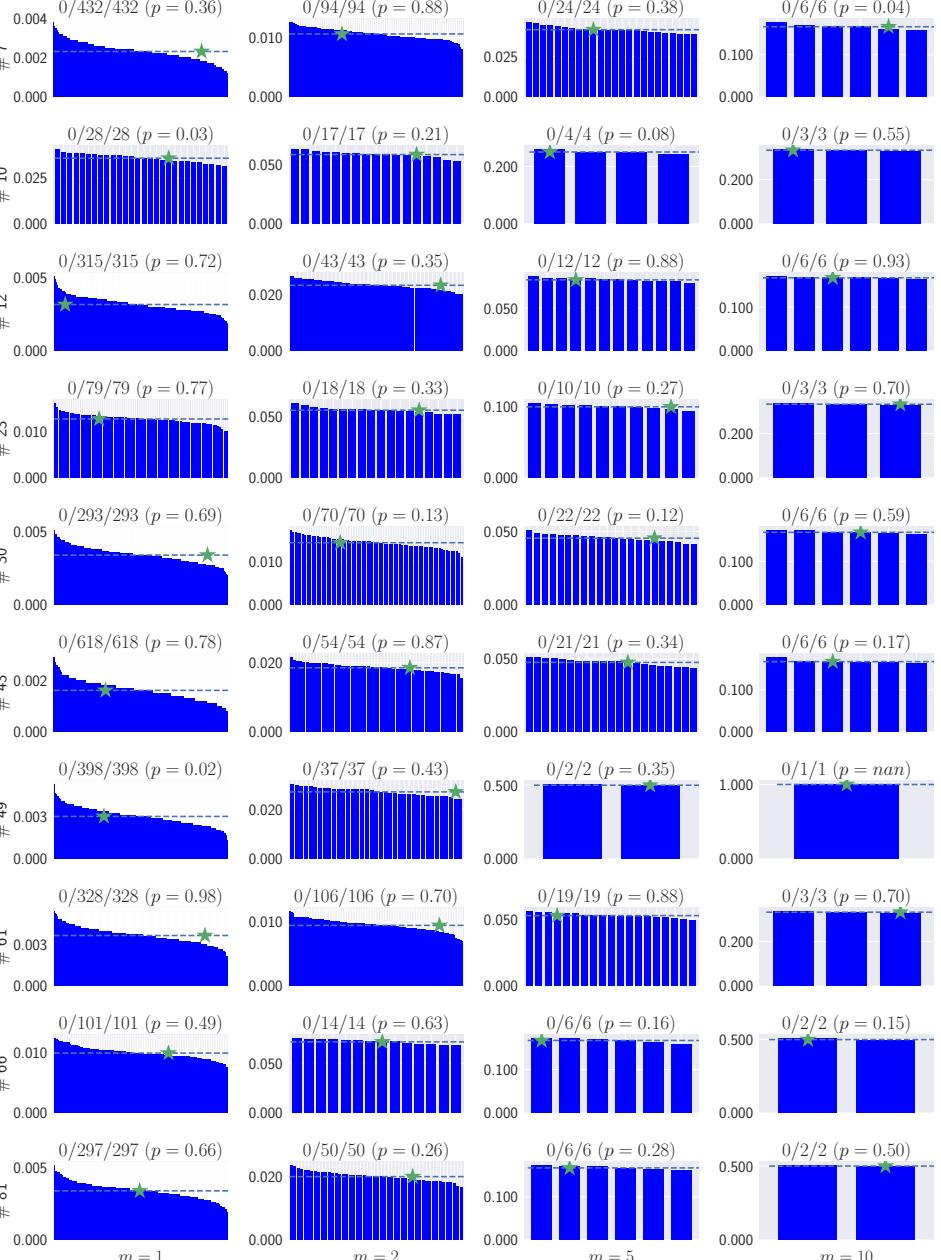
Acknowledgments:

- Experiments were run on NCSA's Blue Waters supercomputer

Canopy



Rejection Sampling



Somatic Mutations Occur at Different Genomic Scales

