

SPhyR: Tumor Phylogeny Estimation from Single-Cell Sequencing Data under Loss and Error

Mohammed El-Kebir – University of Illinois at Urbana Champaign,
Department of Computer Science

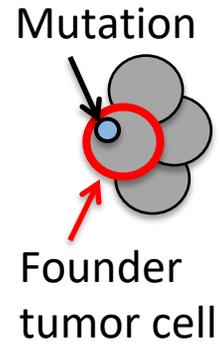
ECCB 2018



Tumorigenesis: (i) Cell Mutation

Clonal Evolution Theory of Cancer

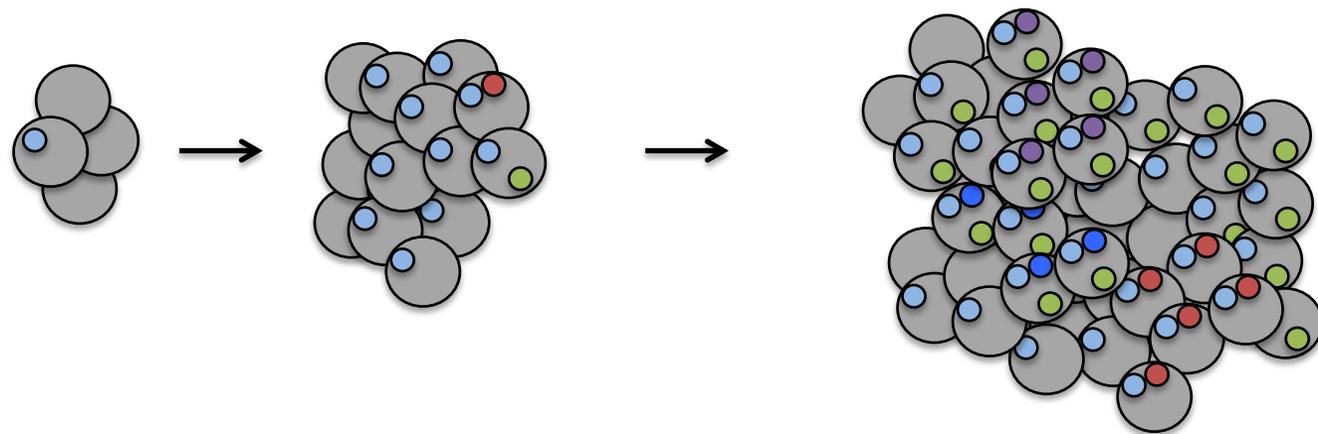
[Nowell, 1976]



Tumorigenesis: (i) Cell Mutation & (ii) Cell Division

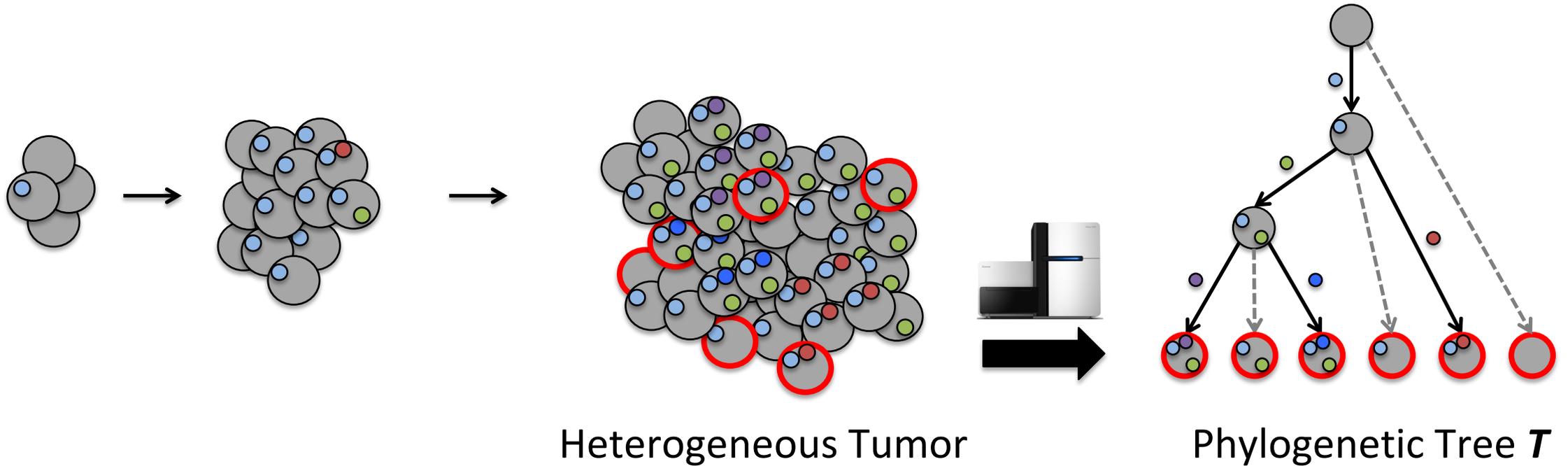
Clonal Evolution Theory of Cancer

[Nowell, 1976]



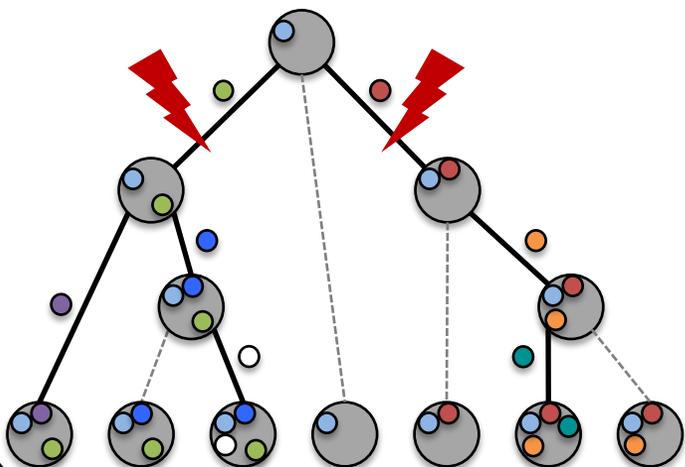
Heterogeneous Tumor

Tumorigenesis: (i) Cell Mutation & (ii) Cell Division

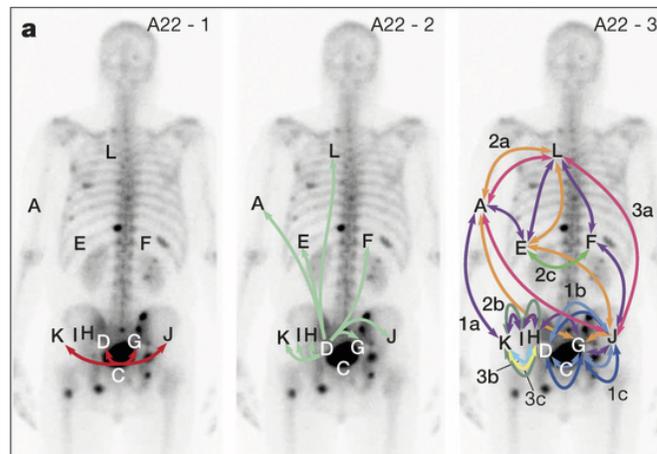


Phylogenies are Key to Understanding Tumorigenesis

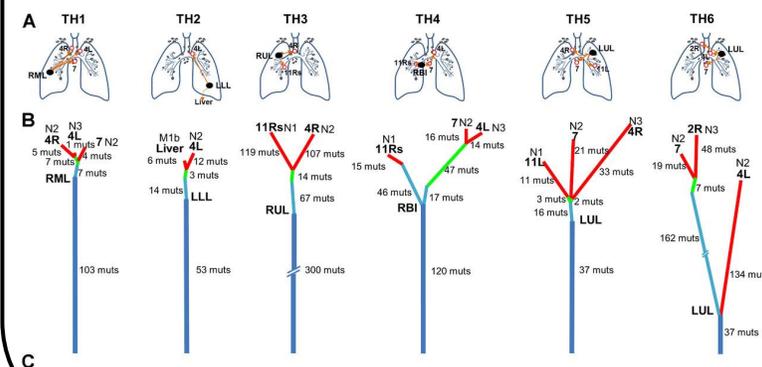
Identify targets for treatment



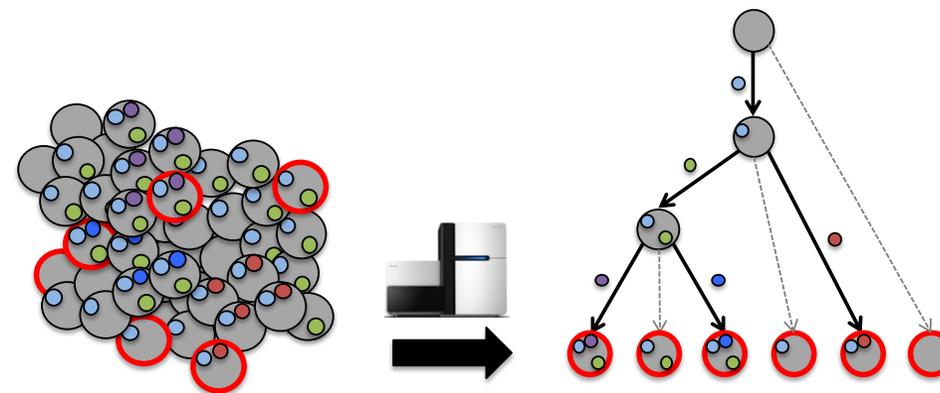
Understand metastatic development



Recognize common patterns of tumor evolution across patients



Goal: Given single-cell DNA sequencing data, find phylogenetic tree T

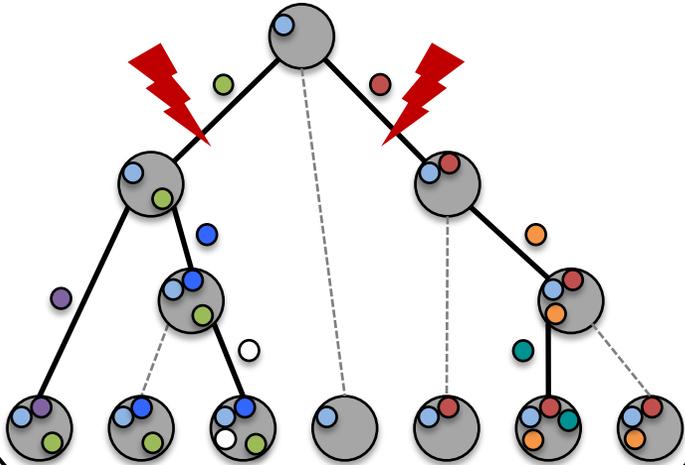


Heterogeneous Tumor

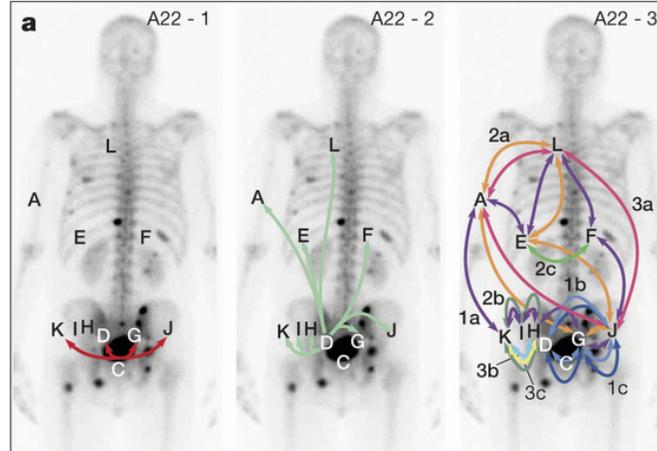
Phylogenetic Tree T

Phylogenies are Key to Understanding Tumorigenesis

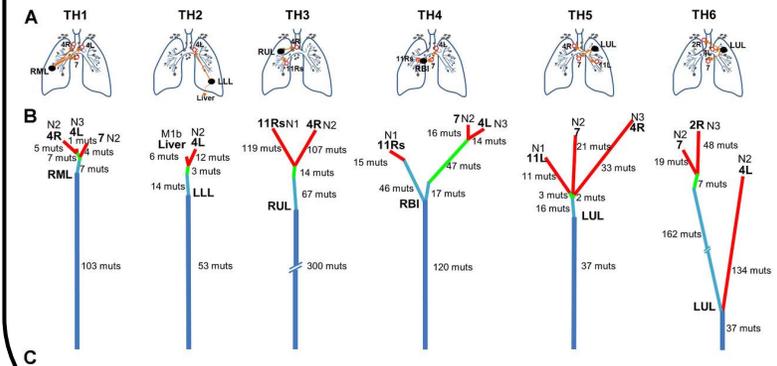
Identify targets for treatment



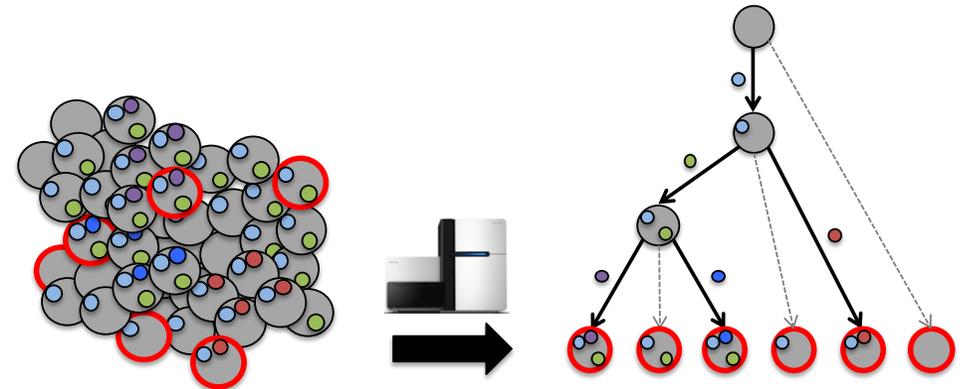
Understand metastatic development



Recognize common patterns of tumor evolution across patients



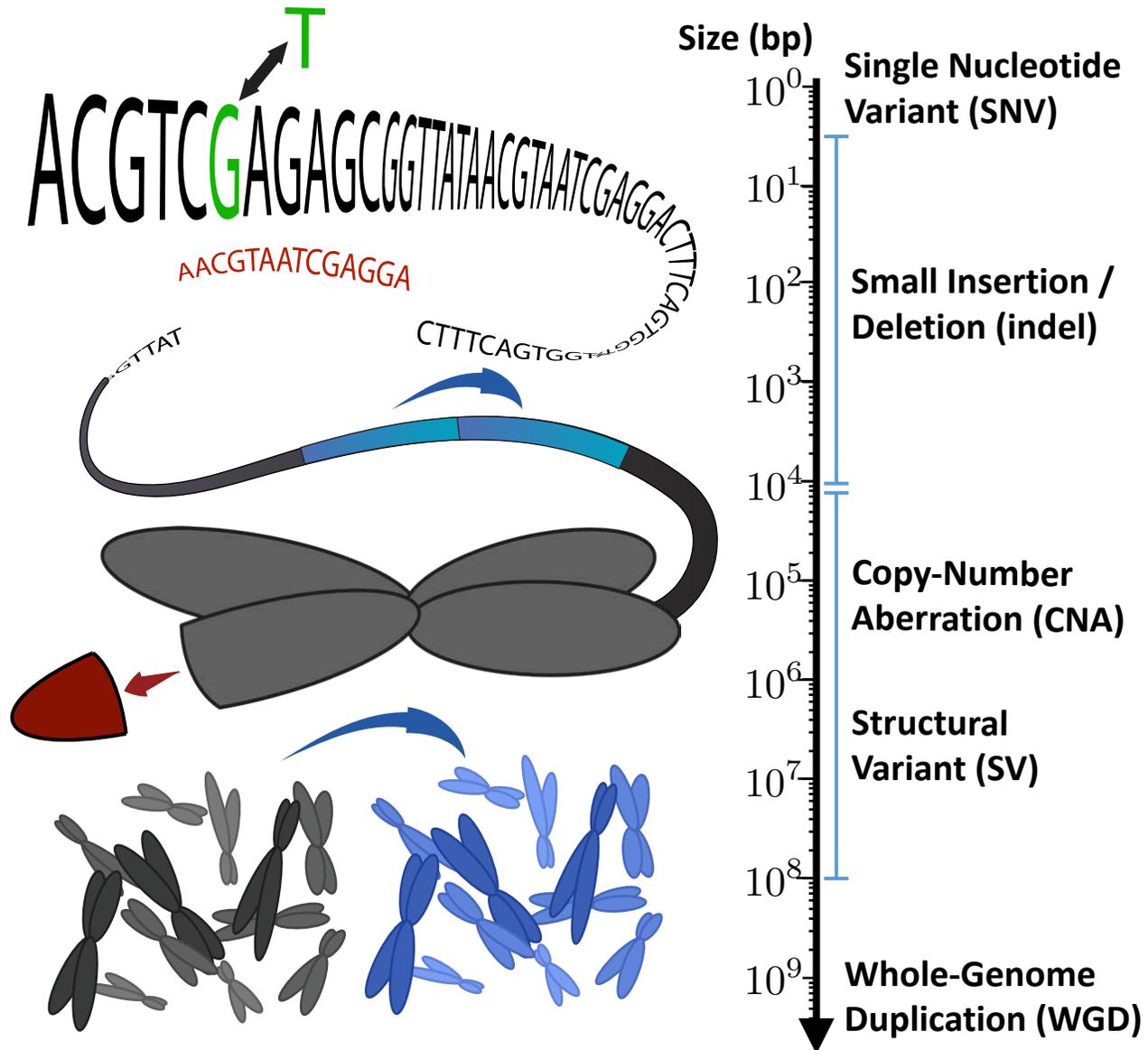
Goal: Given single-cell DNA sequencing data,
find phylogenetic tree T
Requirement: Evolutionary model for somatic
mutations



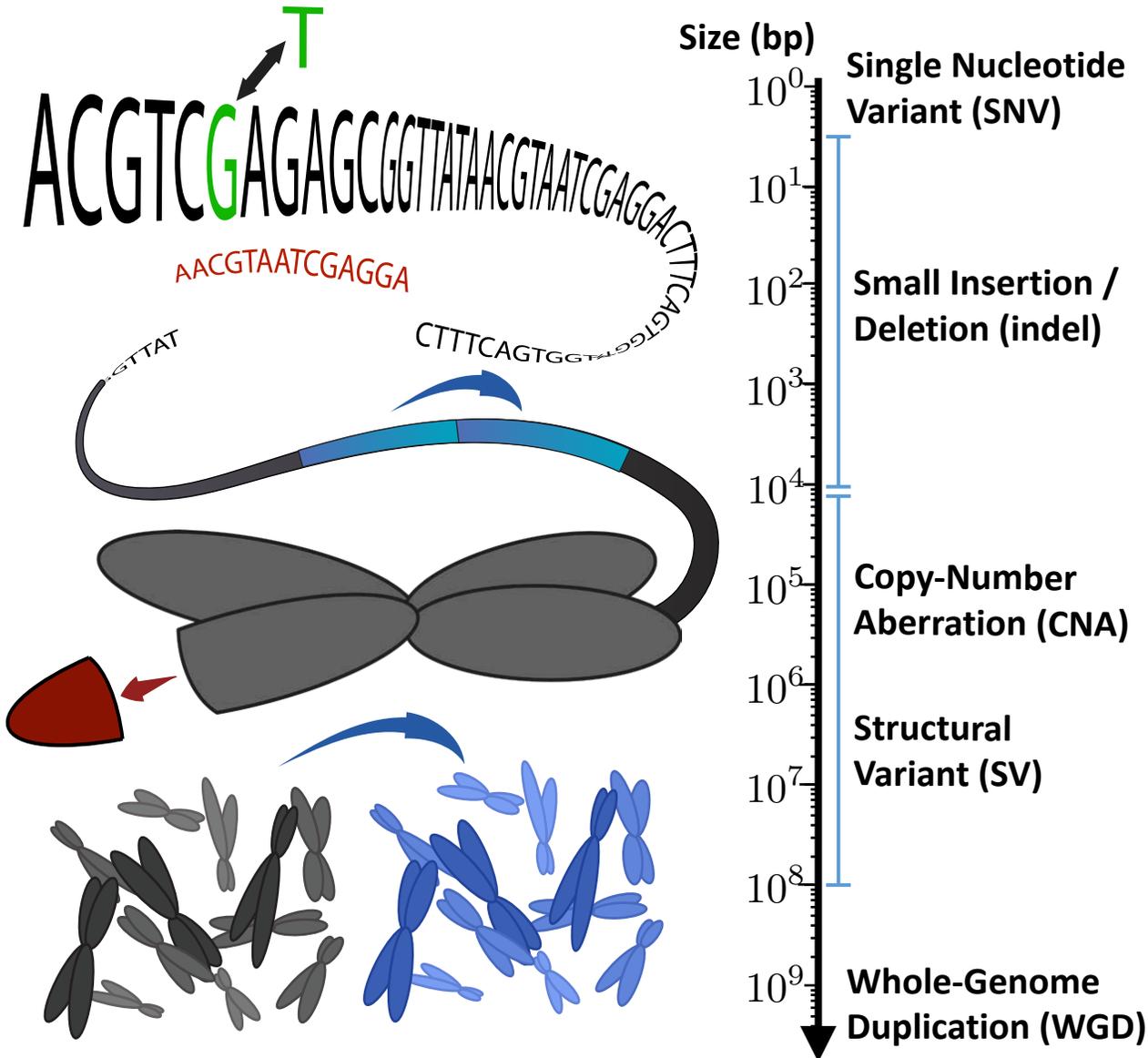
Heterogeneous Tumor

Phylogenetic Tree T

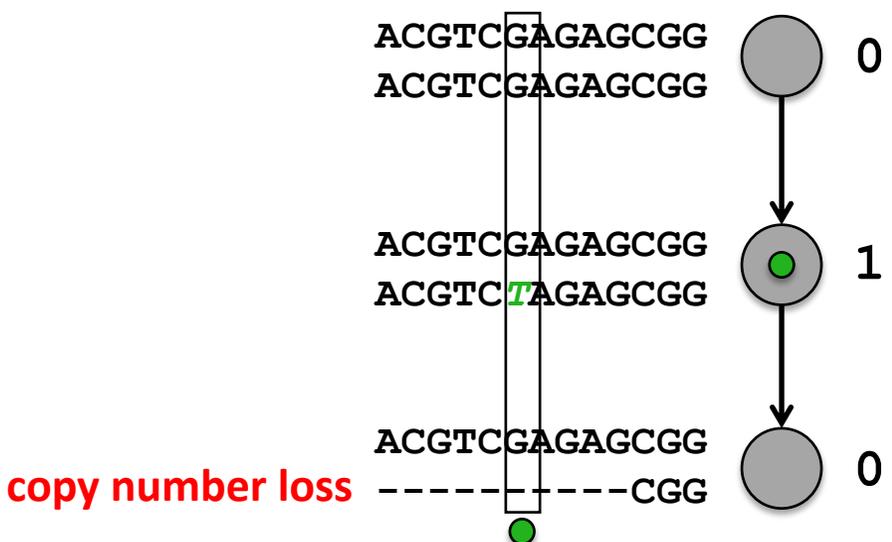
Somatic Mutations Occur at Different Genomic Scales



Infinite Sites Assumption is too Restrictive for SNVs



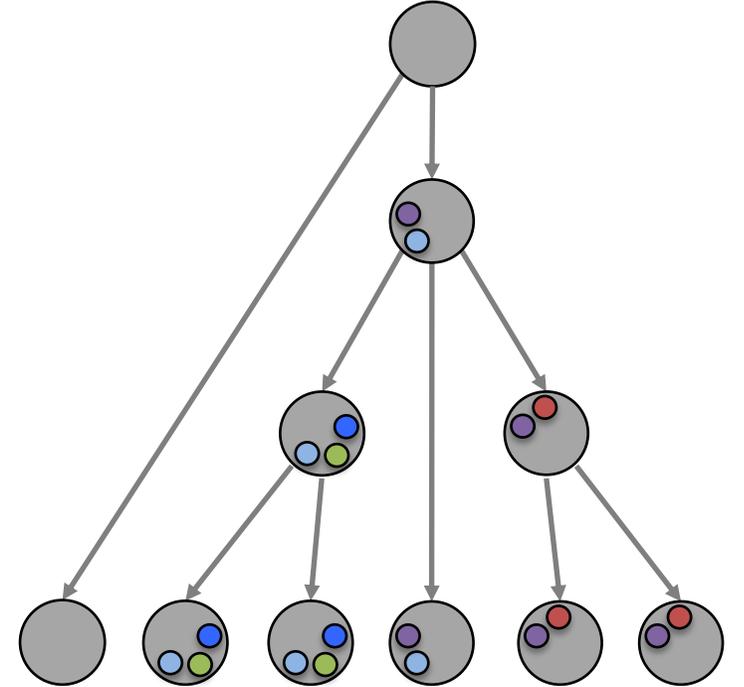
SNVs can be **lost** due to CNAs



- Infinite sites assumption:**
- No parallel evolution of SNVs
 - No loss of SNVs
 - SCITE [Jahn et al. 2016]
 - OncoNEM [Ross and Markowetz, 2016]

Outline

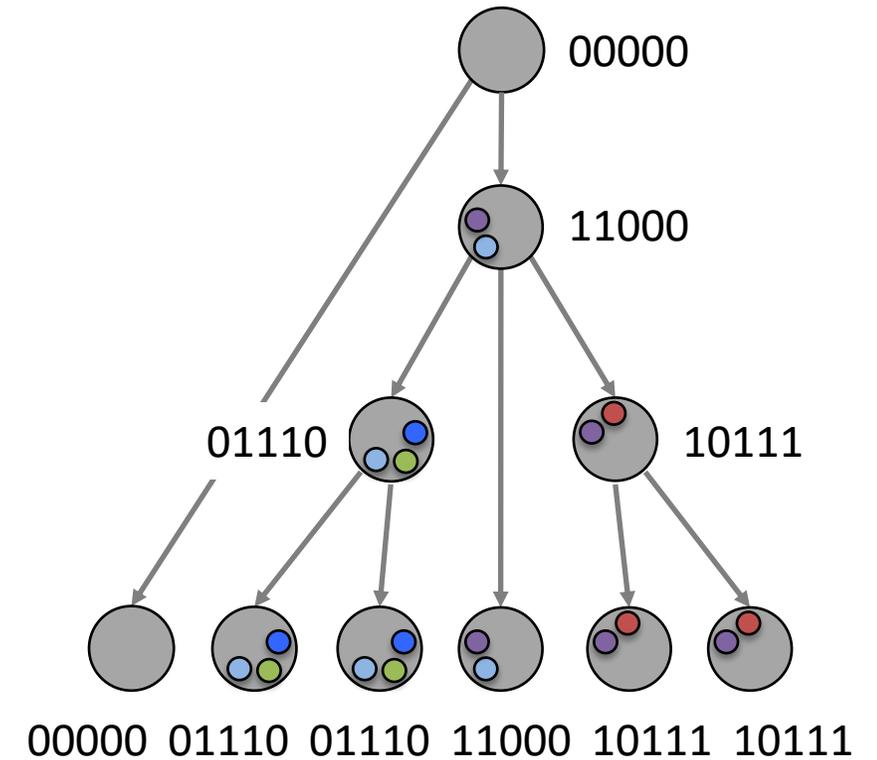
- Perfect data (error free)
 - Problem statement
 - Combinatorial characterization of solutions
 - Exact algorithm
 - Results
- Real data (with errors)
 - Problem statement
 - Heuristic algorithm
 - Results
- Conclusions



k -Dollo Phylogeny (k -DP) Problem

Definition 1. A k -Dollo phylogeny T is a rooted, node-labeled tree subject to the following conditions.

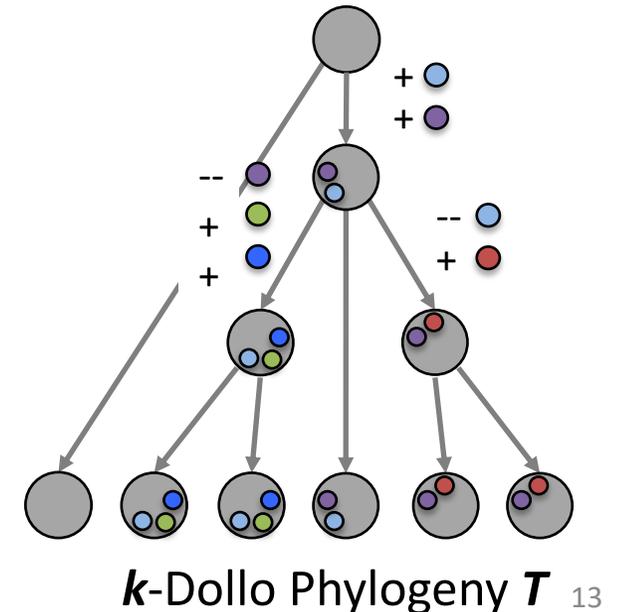
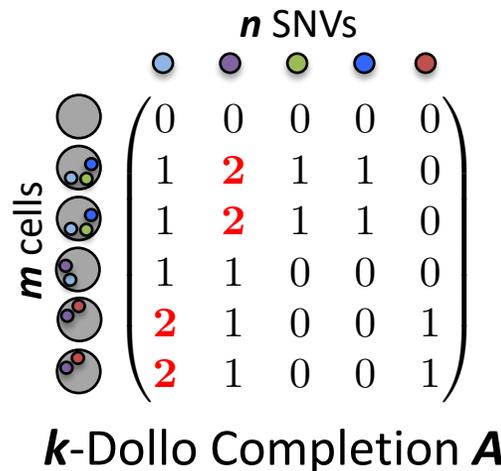
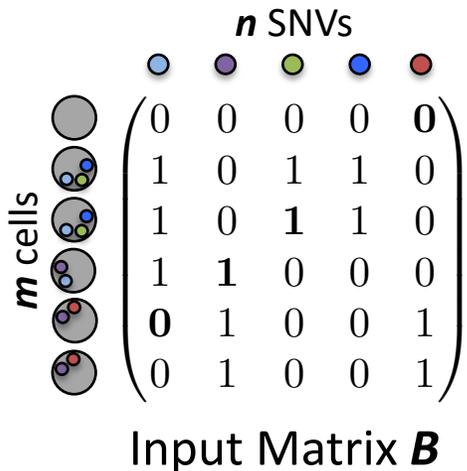
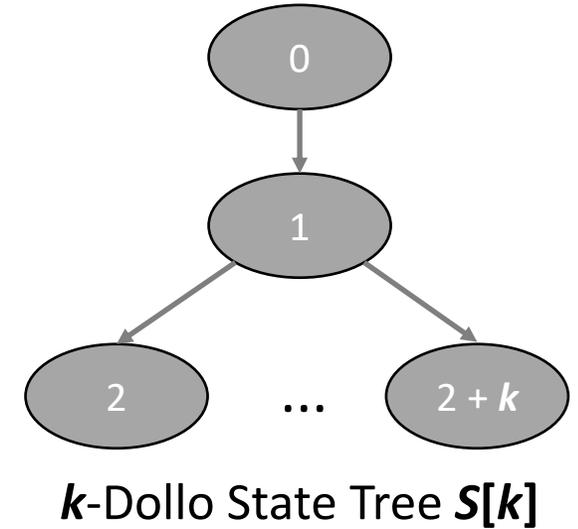
1. Each node v of T is labeled by a vector $\mathbf{b}_v \in \{0, 1\}^n$.
2. The root r of T is labeled by vector $\mathbf{b}_r = [0, \dots, 0]^T$.



Combinatorial Characterization of k -DP

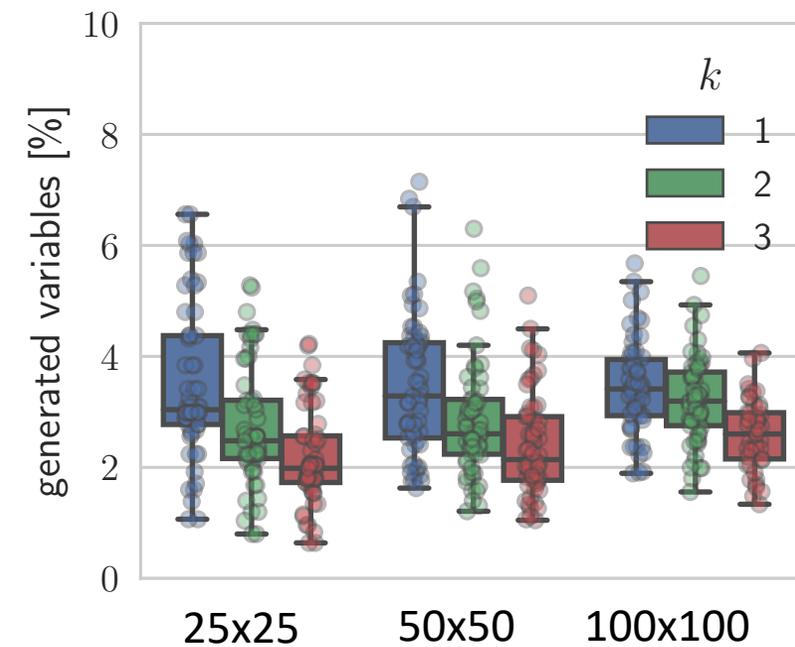
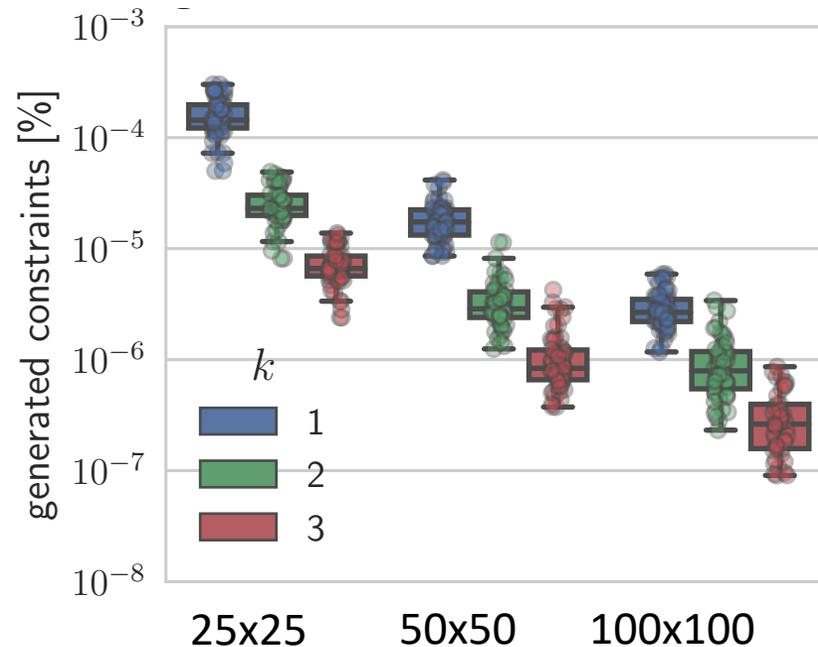
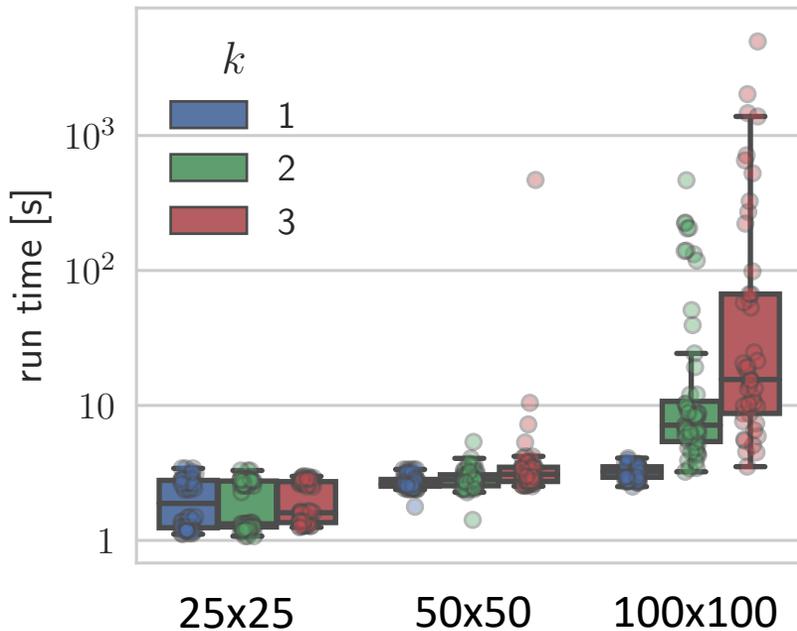
Theorem 3. Let $B \in \{0,1\}^{m \times n}$. The following statements are equivalent.

1. There exists a k -Dollo phylogeny T for B .
2. There exists a k -Dollo completion A of B .
3. There exists a k -completion A of B , and perfect phylogeny T for A whose characters are consistent with $S[k]$.



Results for k -DP

- Naive ILP does not scale and has $O(mnk)$ variables and $O(m^3n^2k^4)$ constraints
- Column and cutting plane generation
 - Introduce variables and constraints only when needed
- Simulations with 60 instances for each m, n and k

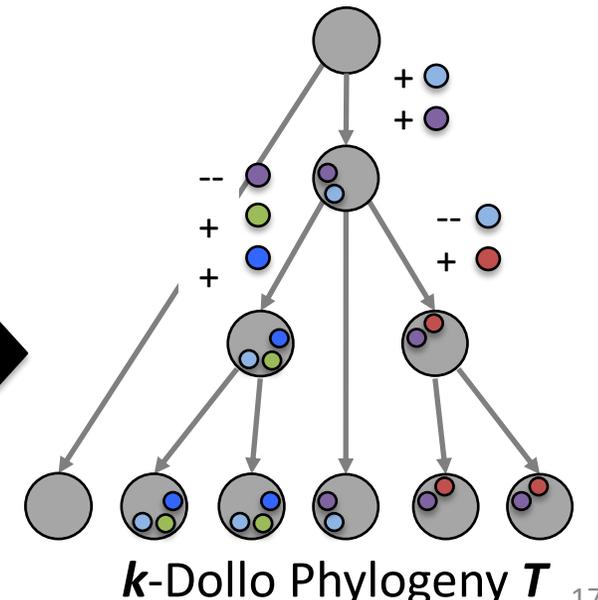
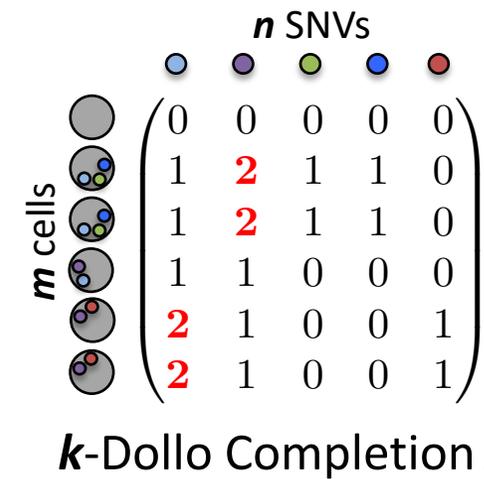
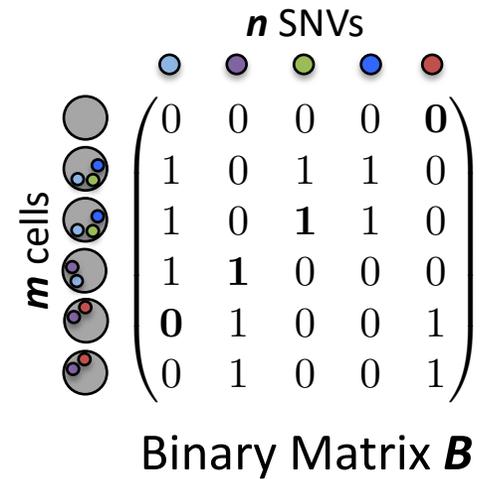
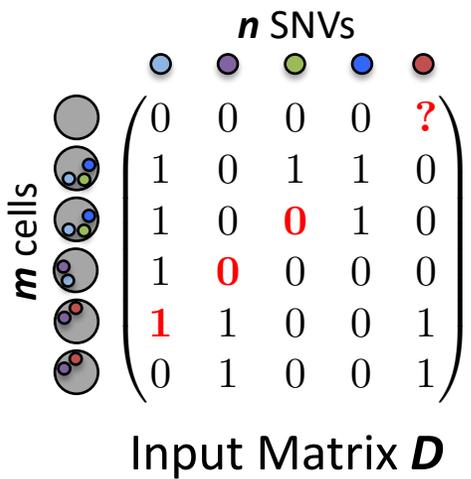


Outline

- Perfect data (error free)
 - Problem statement
 - Combinatorial characterization of solutions
 - Exact algorithm
 - Results
- Real data (with errors)
 - Problem statement
 - Heuristic algorithm
 - Results
- Conclusions

k -Dollo Phylogeny Flip and Cluster (k -DPFC) problem. Given matrix $D \in \{0, 1, ?\}^{m \times n}$, error rates $\alpha, \beta \in [0, 1]$, integers $k, s, t \in \mathbb{N}$, find matrix $B \in \{0, 1\}^{m \times n}$ and tree T such that: (1) B has at most s unique rows and at most t unique columns; (2) $\Pr(D \mid B, \alpha, \beta)$ is maximum; and (3) T is a k -Dollo phylogeny for B .

$$\Pr(D \mid B, \alpha, \beta) = \prod_{p=1}^m \prod_{c=1}^n \begin{cases} \alpha, & d_{p,c} = 1 \text{ and } b_{p,c} = 0 \\ 1 - \alpha, & d_{p,c} = 1 \text{ and } b_{p,c} = 1, \\ \beta, & d_{p,c} = 0 \text{ and } b_{p,c} = 1, \\ 1 - \beta, & d_{p,c} = 0 \text{ and } b_{p,c} = 0, \\ 1, & d_{p,c} = ? \end{cases}$$

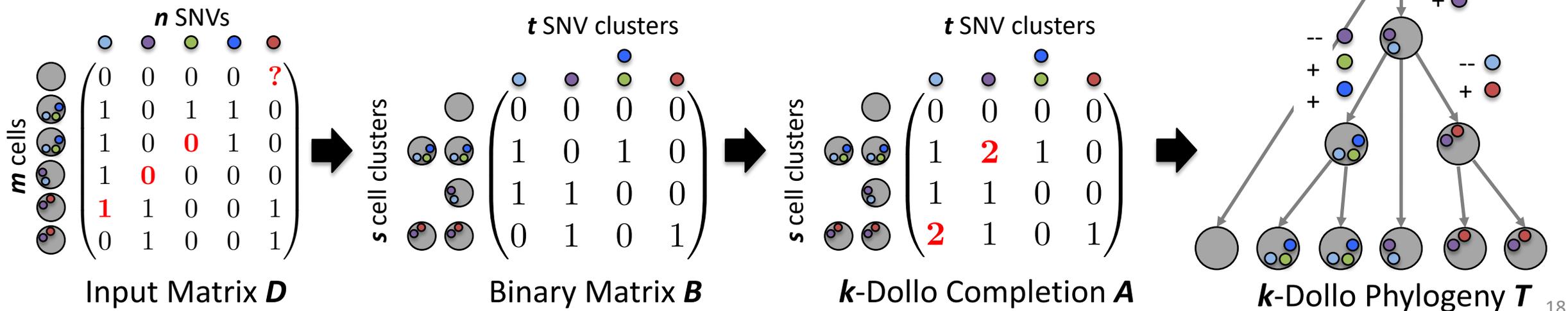


SPhyR: Single-cell Phylogeny Reconstruction

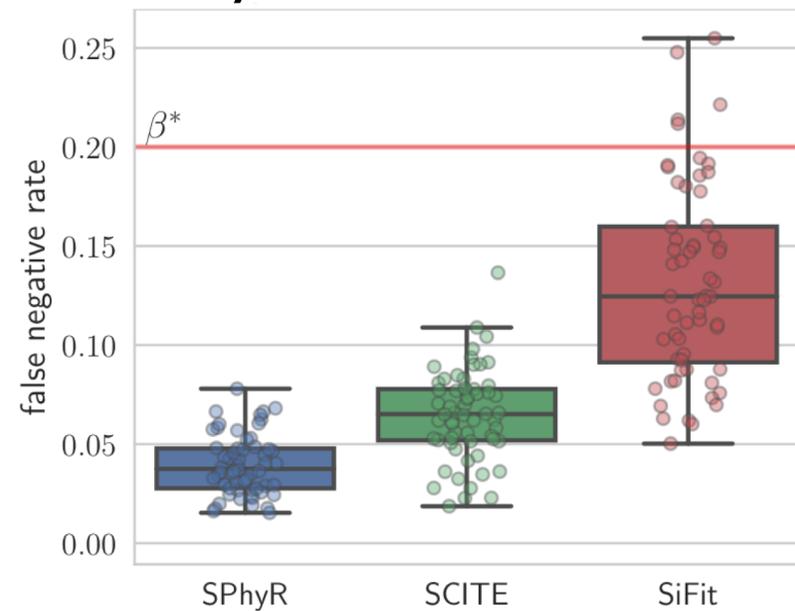
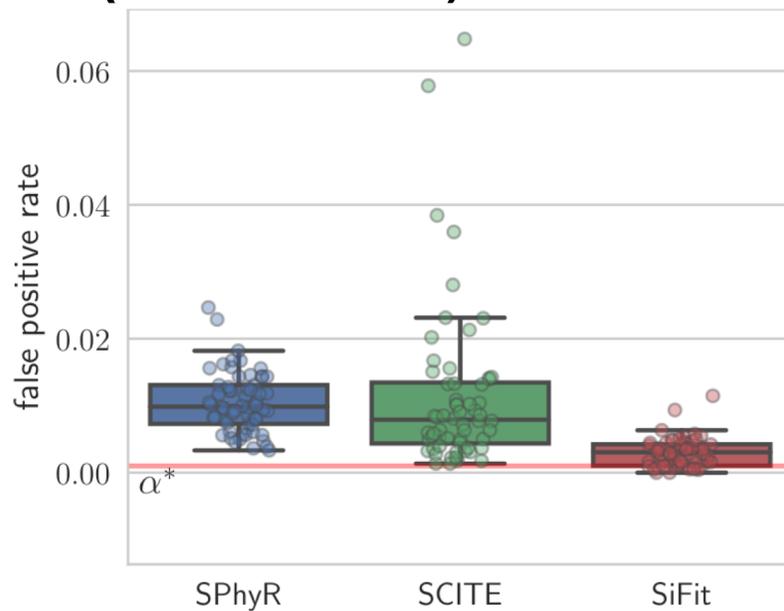
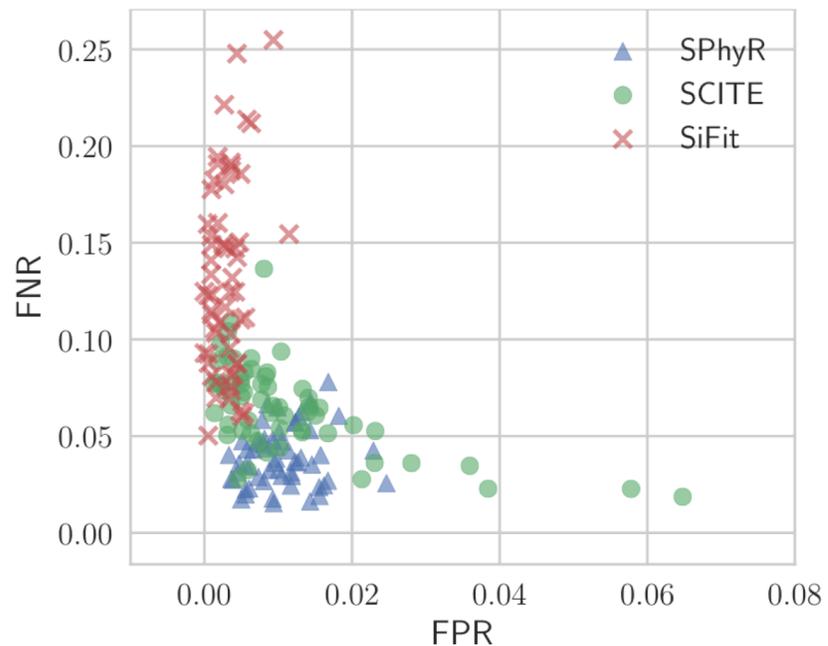
- Coordinate ascent:

1. k-Means with random seed to obtain cell clustering π and SNV clustering ψ
2. ILP to obtain maximum likelihood k -Dollo completion A given D , π and ψ
3. Identify maximum likelihood π given A and ψ
4. Identify maximum likelihood ψ given A and π
5. Repeat until convergence

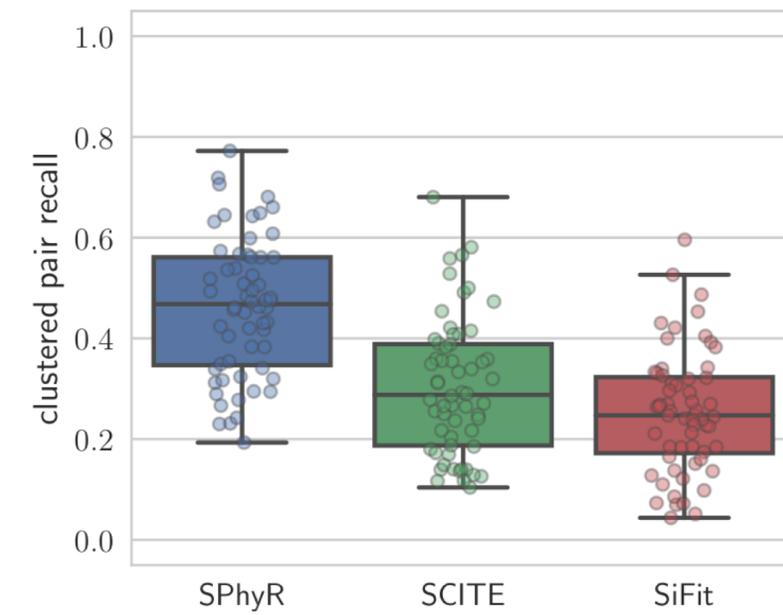
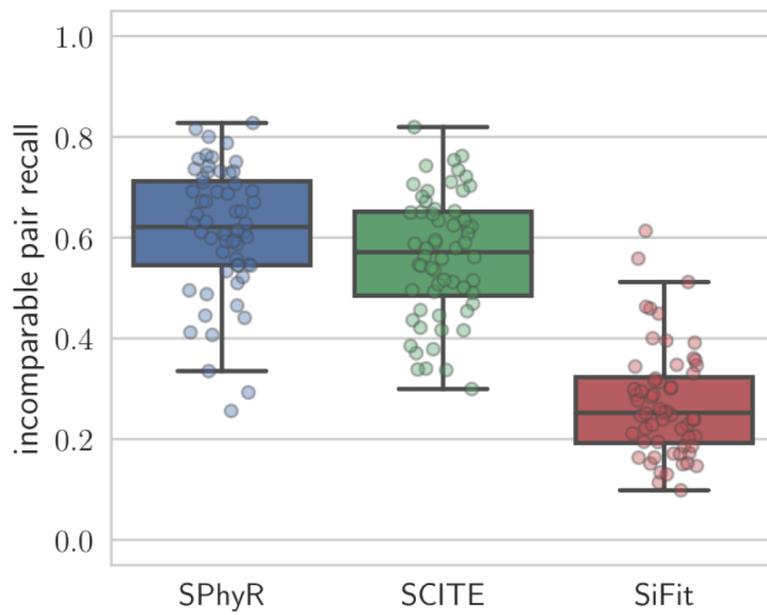
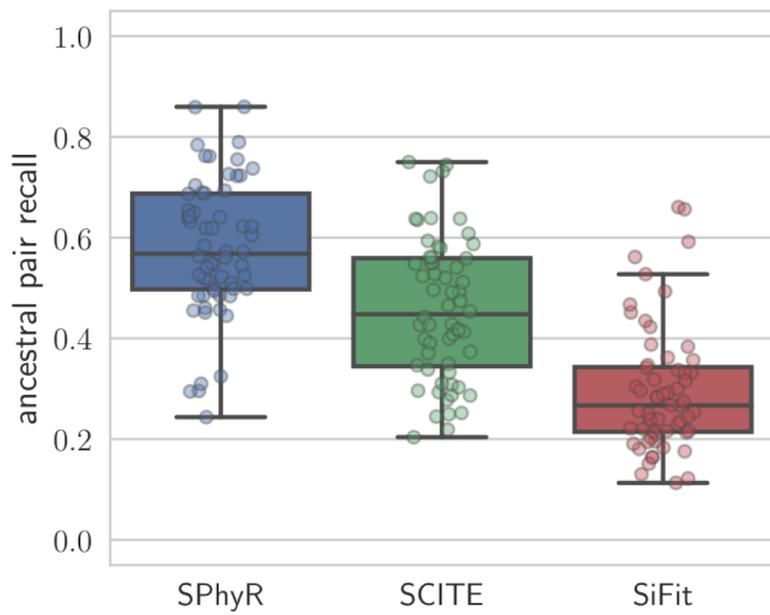
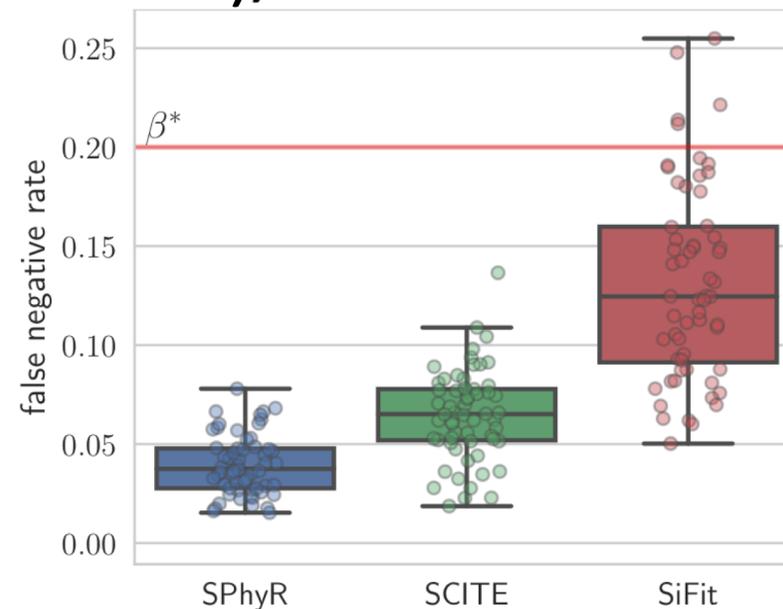
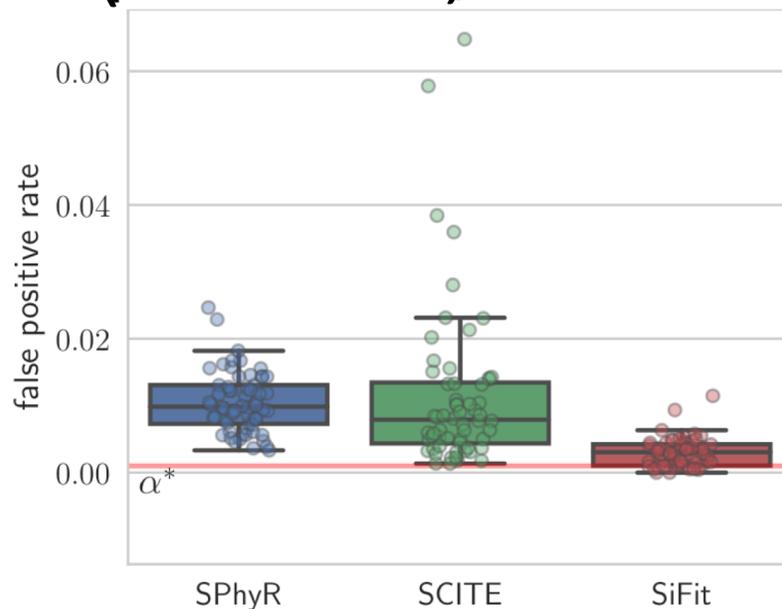
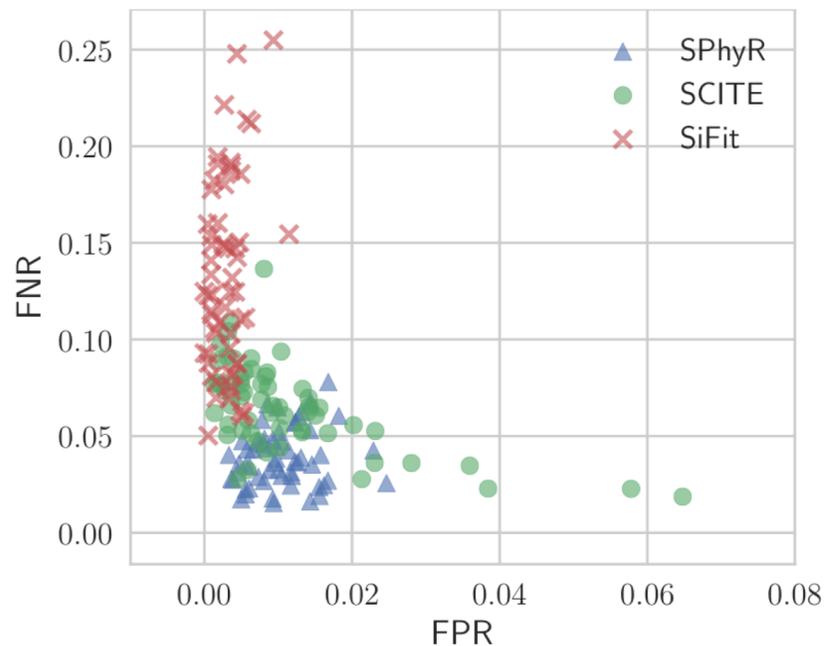
- Available on Github: <https://github.com/elkebir-group/SPhyR>



Simulation Results ($m = 50, n = 50, k = 1$)



Simulation Results ($m = 50, n = 50, k = 1$),



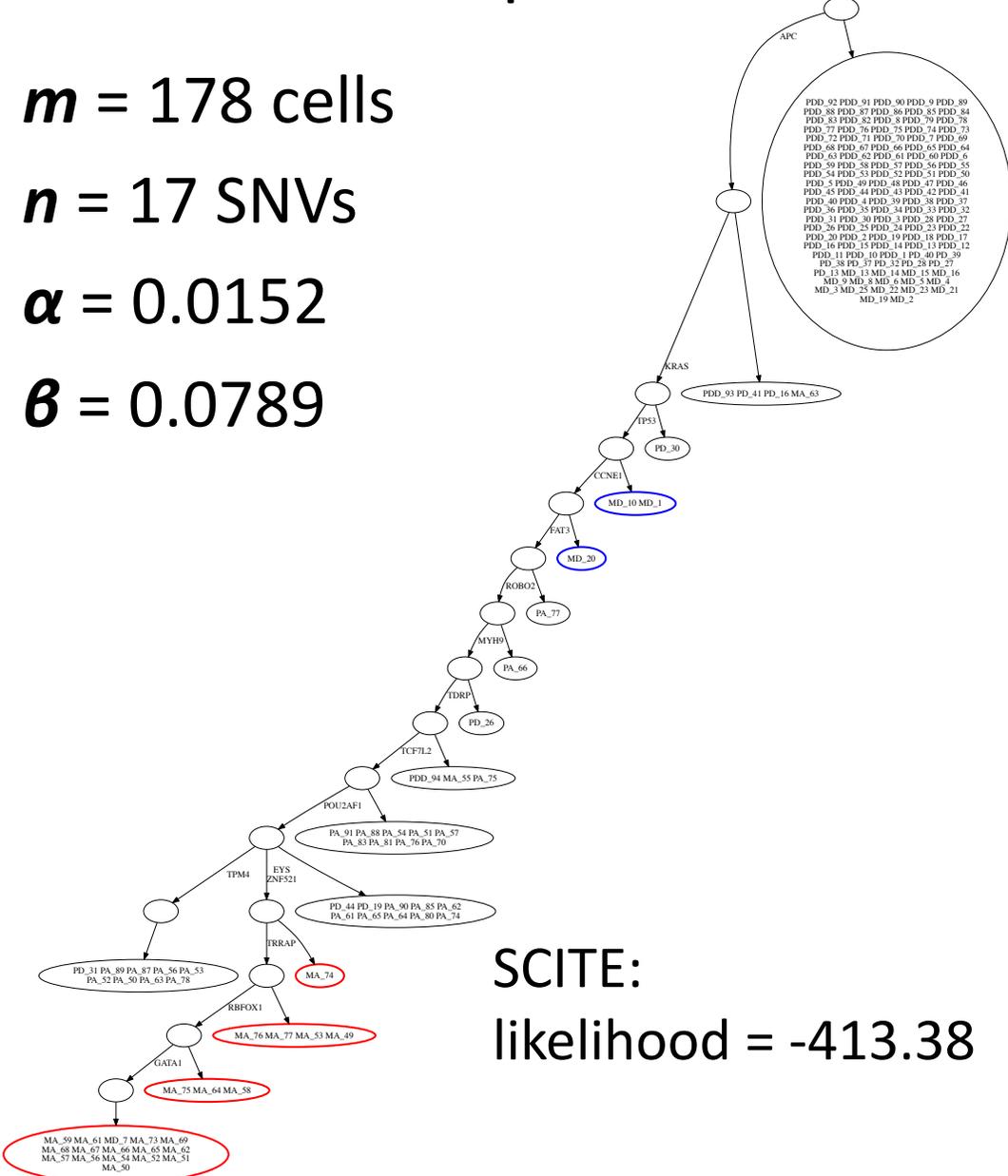
Colorectal patient CRC1 [Leung et al., 2017]

$m = 178$ cells

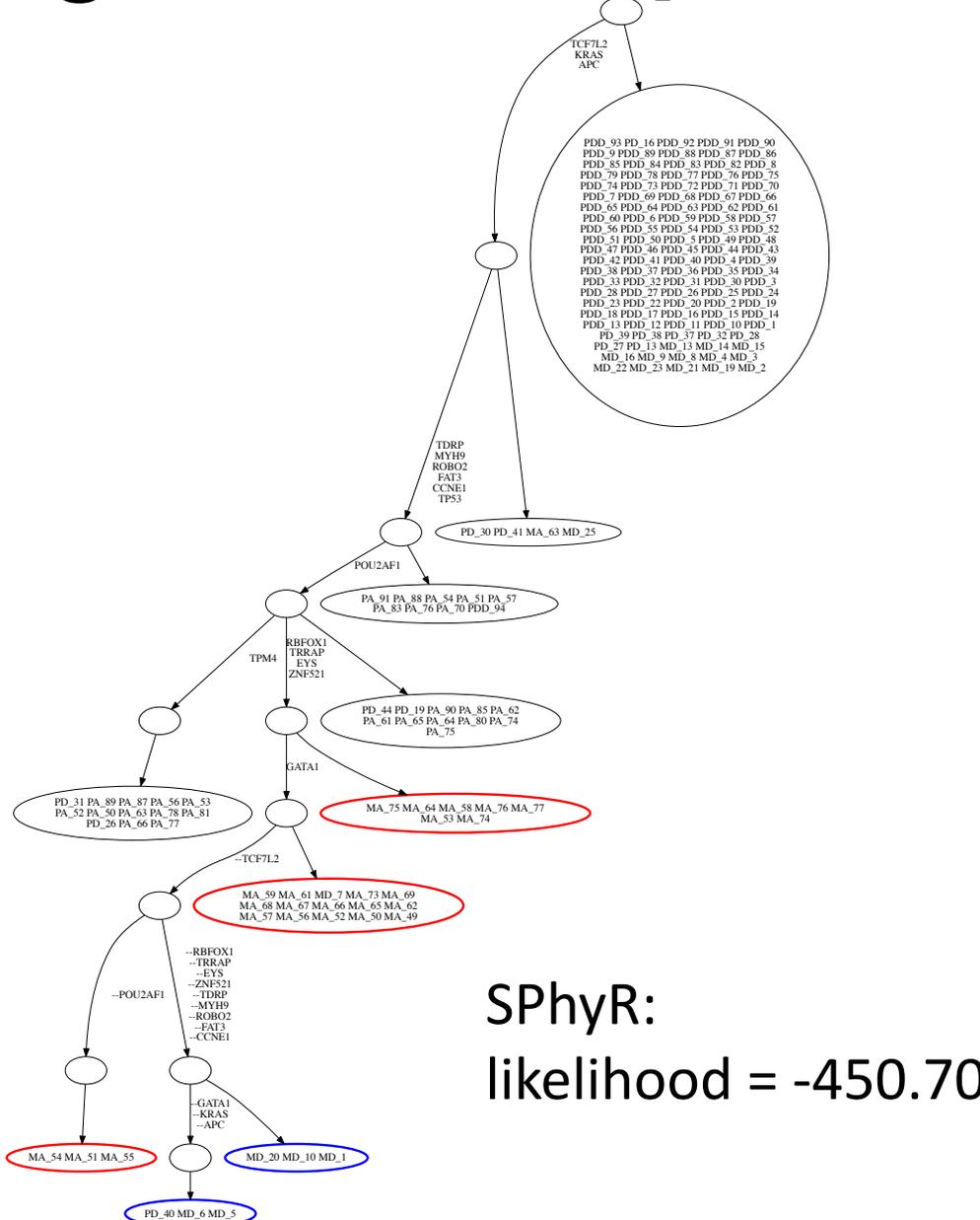
$n = 17$ SNVs

$\alpha = 0.0152$

$\beta = 0.0789$



SCITE:
likelihood = -413.38



SPhyR:
likelihood = -450.70

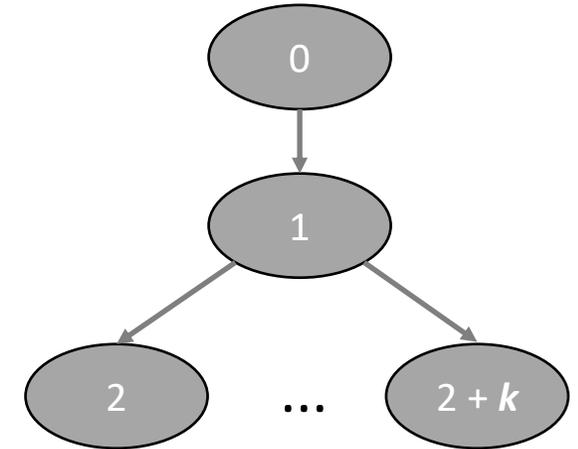
Acknowledgments

- Experiments were run on NCSA Blue Waters

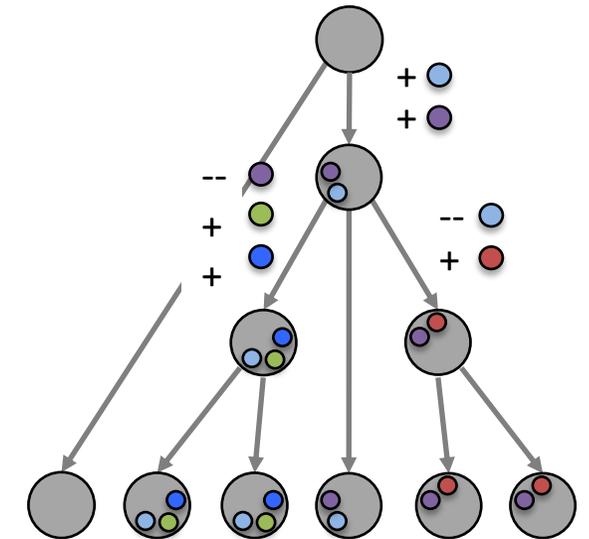
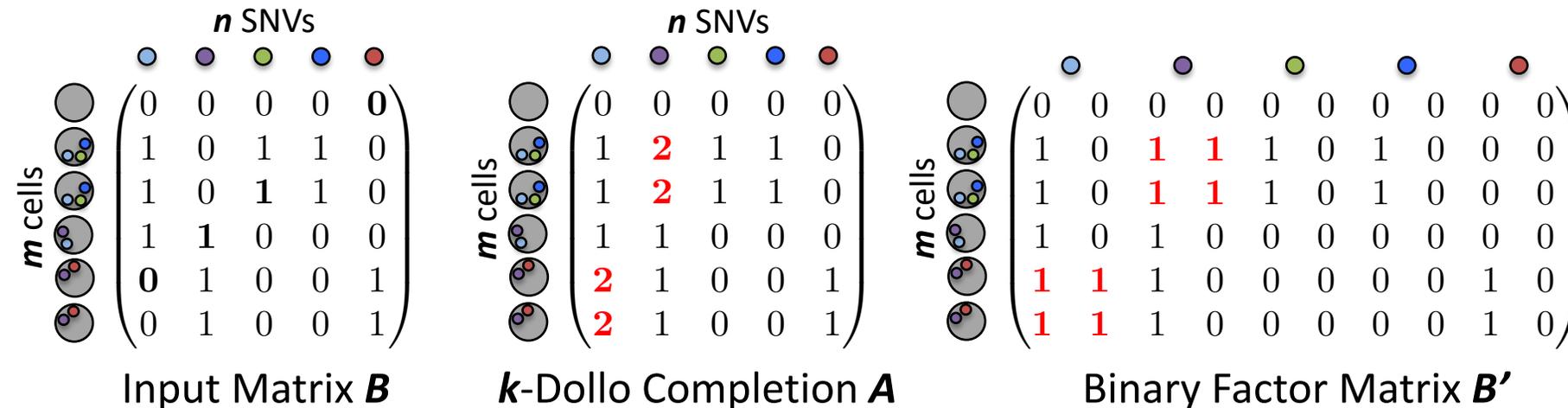
Combinatorial Characterization of k -DP

Theorem 3. Let $B \in \{0,1\}^{m \times n}$. The following statements are equivalent.

1. There exists a k -Dollo phylogeny T for B .
2. There exists a k -Dollo completion A of B .
3. There exists a k -completion A of B such that the binary factor matrix B' of $(A, \mathcal{S}[k])$ is a perfect phylogeny matrix.
4. There exists a k -completion A of B , and perfect phylogeny T for A whose characters are consistent with $\mathcal{S}[k]$.



k -Dollo State Tree $\mathcal{S}[k]$



k -Dollo Phylogeny T 25

