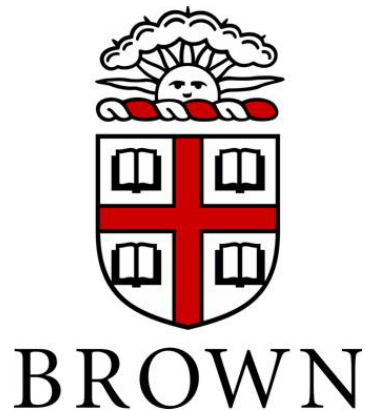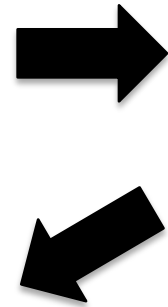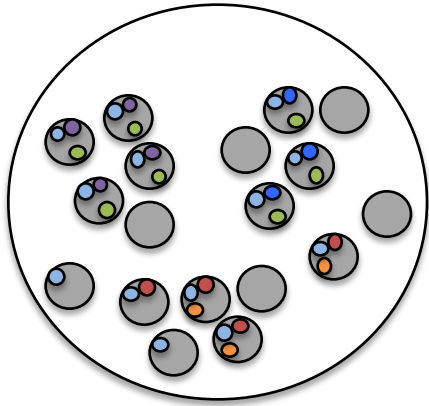# Multi-State Perfect Phylogeny Mixture Deconvolution and Applications to Cancer Sequencing

Mohammed El-Kebir

BROWN

CCMB

# Tumor Evolution as a Two-State Perfect Phylogeny
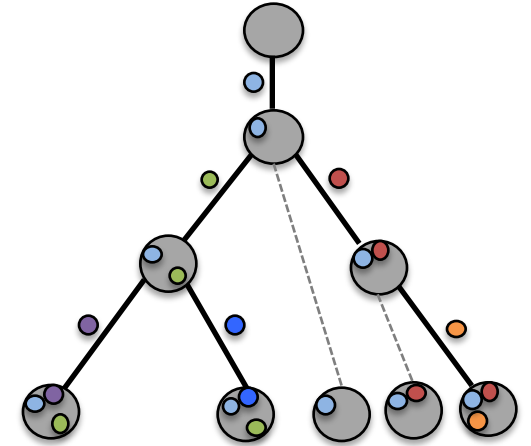
## Tumor snapshot



$O(mn)$

**Single-cell** sequencing

## Two-State Perfect Phylogeny Tree $T$



**Given:**

SNVs

$$M = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}$$

leaves of $T$

**Find:**

Two-state perfect phylogeny tree $T$

**Assumptions:**
- No copy number aberrations
- Infinite sites assumption

**States:**

0 : non-mutated

1 : mutated

| Seq. method | Mixing | Inferring $T$ |
|---|---|---|
| single-cell | no | two-state perfect phylogeny [Gusfield, 1991] |

2

# Tumor Evolution as a Two-State Perfect Phylogeny



**Tumor snapshot**

$S_3$   $S_2$   $S_1$

**NP-complete**

**Bulk** sequencing

**Two-State Perfect Phylogeny Tree $T$**

$U$   0.8 $S_3$   $S_2$   0.6   0.2   0.2   0.4 $S_1$

**Given:**

mutations

VAFs $F = \begin{bmatrix} 0.4 & 0.0 & 0.0 & 0.0 & 0.3 & 0.2 \\ 0.3 & 0.3 & 0.0 & 0.3 & 0.0 & 0.0 \\ 0.4 & 0.4 & 0.4 & 0.0 & 0.0 & 0.0 \end{bmatrix} \begin{matrix} s_1 \\ s_2 \\ s_3 \end{matrix}$ samples

**Find:**

Two-state perfect phylogeny tree $T$
Mixing proportions $U$

**Variant Allele Frequency (VAF)**:
*Fraction* of reads covering position of single-nucleotide variant (SNV) that contain *variant* allele

GCG**G**ACGT
C**G**ACGTGA
GCG**T**ACGT
**T**ACGTGGA
TGCG**G**ACG
...GAGAAAGCTGCG**G**ACGTGGACGA...

VAF = 2/5 = 0.4

| Seq. method | Mixing | Inferring $T$ |
|---|---|---|
| single-cell | no | two-state perfect phylogeny [Gusfield, 1991] |
| bulk | yes | TrAp [Strino *et al.*, 2013] Rec-BTP [Hajirasouliha *et al.*, 2014] PhyloSub [Jiao *et al.*, 2014] Clomial [Zare *et al.*, 2014] Binary $F$ [Hajirasouliha *et al.*, 2014] CITUP [Malikic *et al.*, 2015] BitPhylogeny [Yuan *et al.*, 2015] LICHeE [Popic *et al.*, 2015] AncesTree [El-Kebir, Oesper *et al.*, 2015] ... |

2

# Tumor Evolution as a Two-State Perfect Phylogeny

**Tumor snapshot**

S₃
S₂
S₁

**Bulk** sequencing

**NP-complete**

**Two-State Perfect Phylogeny Tree $T$**

$U$   0.8  $S_3$   $S_2$  0.6   0.2  0.2  0.4  $S_1$

**Given:**

mutations

$$\text{VAFs } \mathbf{F} = \begin{bmatrix} 0.4 & 0.0 & 0.0 & 0.0 & 0.3 & 0.2 \\ 0.3 & 0.3 & 0.0 & 0.3 & 0.0 & 0.0 \\ 0.4 & 0.4 & 0.4 & 0.0 & 0.0 & 0.0 \end{bmatrix} \begin{matrix} s_1 \\ s_2 \\ s_3 \end{matrix}$$

samples

**Find:**

Two-state perfect phylogeny tree $T$
Mixing proportions $U$

**States:**      rescale VAFs to CCFs

0 : non-mutated

1 : mutated

2 : CN loss-of-heterozygosity

3 : amplification

…                                              ….
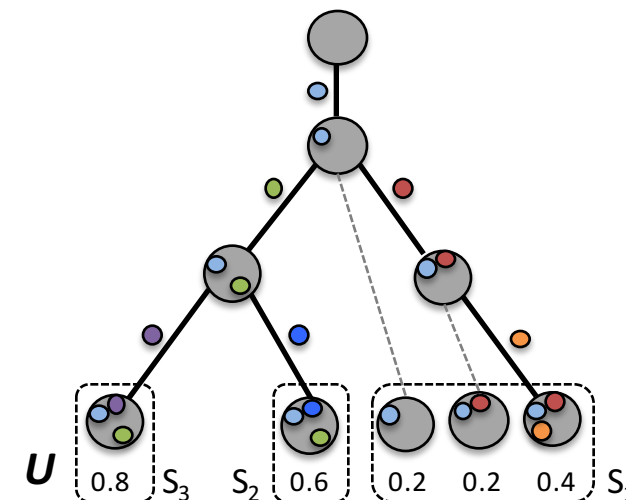
| Seq. method | Mixing | Inferring $T$ |
|---|---|---|
| single-cell | no | two-state perfect phylogeny [Gusfield, 1991] |
| bulk | yes | TrAp [Strino *et al.,* 2013] Rec-BTP [Hajirasouliha *et al.,* 2014] PhyloSub [Jiao *et al.,* 2014] Clomial [Zare *et al.,* 2014] Binary $F$ [Hajirasouliha *et al.,* 2014] CITUP [Malikic *et al.,* 2015] BitPhylogeny [Yuan *et al.,* 2015] LICHeE [Popic *et al.,* 2015] AncesTree [El-Kebir, Oesper *et al.,* 2015] … |

2

# Tumor Evolution as a <span style="color:red">Multi</span>-State Phylogeny

## Tumor snapshot



## Two-State Perfect Phylogeny Tree **T**



$U$    0.8   $S_3$    $S_2$   0.6    0.2   0.2   0.4   $S_1$

**NP-complete**

**Bulk** sequencing

## Given:

mutations

VAFs $F$ = $\begin{bmatrix} 0.4 & 0.0 & 0.0 & 0.0 & 0.3 & 0.2 \\ 0.3 & 0.3 & 0.0 & 0.3 & 0.0 & 0.0 \\ 0.4 & 0.4 & 0.4 & 0.0 & 0.0 & 0.0 \end{bmatrix}$ $\begin{matrix} s_1 \\ s_2 \\ s_3 \end{matrix}$ samples

## Find:

Two-state perfect phylogeny tree **T**
Mixing proportions **U**

## States:  more than > 2 states

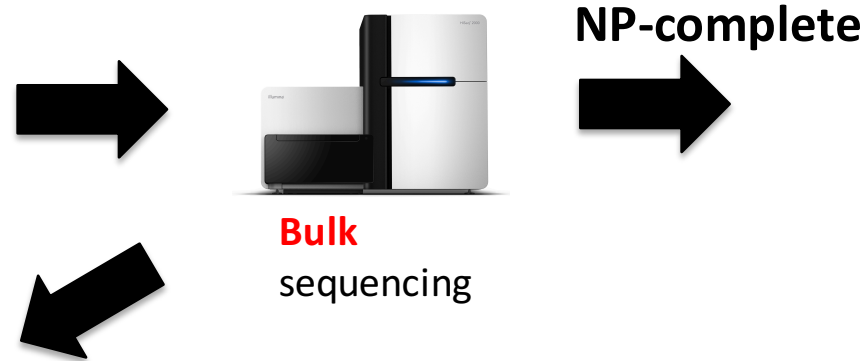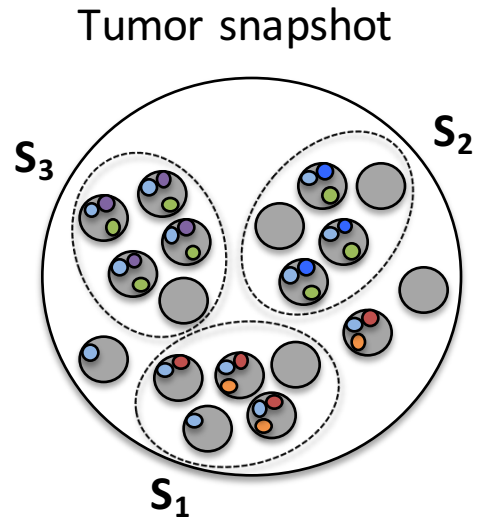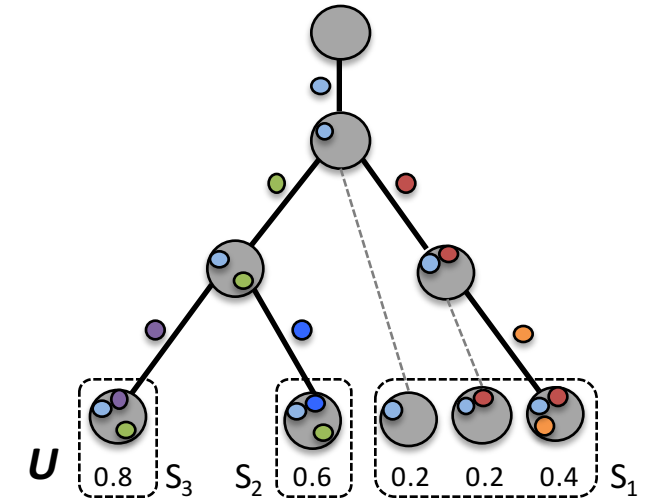0 : non-mutated

1 : mutated

2 : CN loss-of-heterozygosity

3 : amplification

…

| Seq. method | Mixing | Inferring *T* |
|---|---|---|
| single-cell | no | two-state perfect phylogeny [Gusfield, 1991] |
| bulk | yes | TrAp [Strino *et al.*, 2013] Rec-BTP [Hajirasouliha *et al.*, 2014] PhyloSub [Jiao *et al.*, 2014] Clomial [Zare *et al.*, 2014] Binary *F* [Hajirasouliha *et al.*, 2014] CITUP [Malikic *et al.*, 2015] BitPhylogeny [Yuan *et al.*, 2015] LICHeE [Popic *et al.*, 2015] AncesTree [El-Kebir, Oesper *et al.*, 2015] … |

2

# Outline

• Problem Statement

• Combinatorial Characterization of Solutions

• Application to Cancer Sequencing

# Problem Statement

**Two-State Perfect Phylogeny:**

Infinite sites assumption: a character changes state once

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

*n* mutations

# Problem Statement

**Two-State Perfect Phylogeny:**

Infinite sites assumption: a character changes state once



***n* mutations**

$$F = \begin{pmatrix} 0.8 & 0.8 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.6 & 0.6 & 0.0 & 0.6 & 0.0 & 0.0 \\ 0.8 & 0.0 & 0.0 & 0.0 & 0.6 & 0.4 \end{pmatrix}$$ *m* samples

taxa

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

Usage Matrix ***U***

mutations

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

Complete Two-State
Perfect Phylogeny ***B* / *T***

***U*=[$u_{pj}$]** is a ***usage matrix*** iff

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

1-1

$S_1$   $S_2$   $S_3$

0.8   0.6   0.2   0.2   0.4

# Problem Statement

**Two-State Perfect Phylogeny:**

Infinite sites assumption: a character changes state once

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

$n$ mutations / $m$ samples

$$= \begin{pmatrix} 0.0 & 0.0 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.6 & 0.0 & 0.0 \\ 0.2 & 0.0 & 0.0 & 0.0 & 0.2 & 0.4 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 & 1 \end{pmatrix}$$

taxa / samples

Usage Matrix $U$

mutations / taxa

Complete Two-State Perfect Phylogeny $B / T$

1-1

$U=[u_{pj}]$ is a *usage matrix* iff

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

**VAF Factorization Problem (VAFFP):** [El-Kebir, Oesper et al., 2015]

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

$S_1$    $S_2$    $S_3$

0.8    0.6    0.2   0.2   0.4

# Problem Statement

## Two-State Perfect Phylogeny:

Infinite sites assumption: a character changes state once

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

*n* mutations

*m* samples

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

taxa

samples

Usage Matrix **U**

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

mutations

taxa

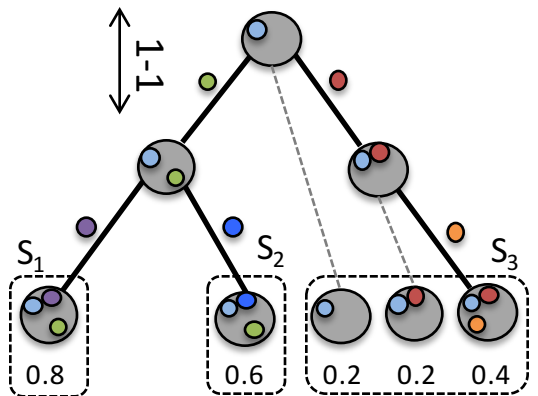Complete Two-State Perfect Phylogeny **B / T**

**U**=[$u_{pj}$] is a ***usage matrix*** iff

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

**VAF Factorization Problem (VAFFP):** [El-Kebir, Oesper et al., 2015]

Given **F**, find **U** and **B** such that **F = U B**

1-1

S₁ 0.8   S₂ 0.6   S₃ 0.2 0.2 0.4

## Multi-State Perfect Phylogeny:

Infinite alleles assumption: a character changes to a state once

characters

$$\begin{pmatrix} 0.1 & 0.8 \\ 0.7 & 0.0 \end{pmatrix}$$
$F_0$

samples

$$\begin{pmatrix} 0.2 & 0.2 \\ 0.3 & 0.4 \end{pmatrix}$$
$F_1$

$$\begin{pmatrix} 0.7 & 0.0 \\ 0.0 & 0.6 \end{pmatrix}$$
$F_2$

*m* samples   *k* states   *n* characters

Frequency Tensor **F**

$$\begin{pmatrix} 1 & 1 \\ 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 0 \end{pmatrix}$$
$A_0$

$$\begin{pmatrix} 0 & 0 \\ 1 & 1 \\ 0 & 0 \\ 0 & 1 \\ 0 & 0 \end{pmatrix}$$
$A_1$

$$\begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix}$$
$A_2$

1-*

[0,0] $V_{(*,0)}$

(**c**,2)   (**d**,1)

[0,1]

(**c**,1)   (**d**,2)

$V_{(c,2)}$   $V_{(*,0)}$   $V_{(c,1)}$   $V_{(d,1)}$   $V_{(d,2)}$
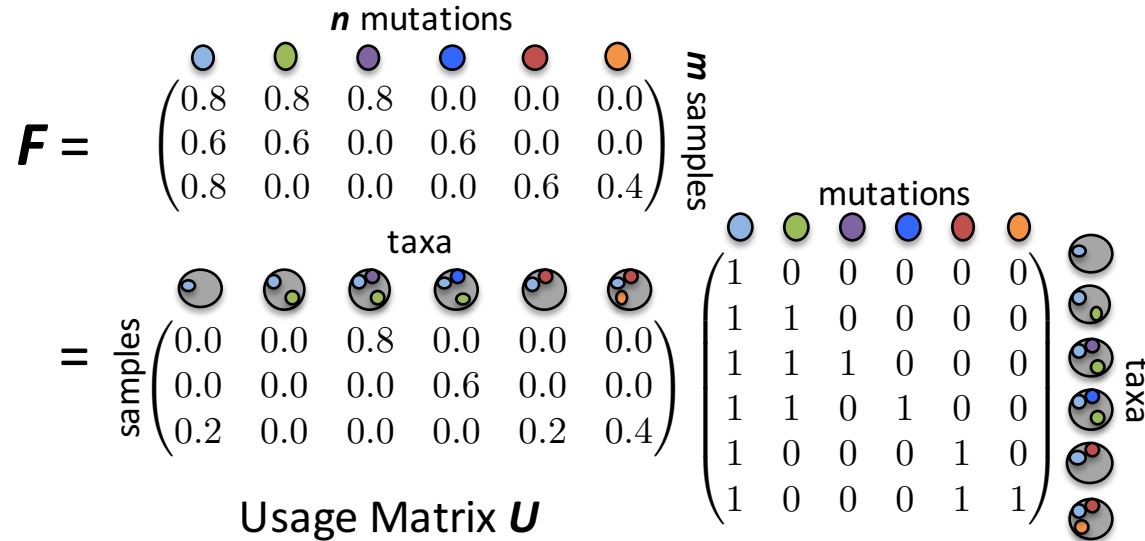
[2,0]   [0,0]   [1,1]   [0,1]   [0,2]

Complete Multi-State Perfect Phylogeny **A / T**

# Problem Statement

## Two-State Perfect Phylogeny:

Infinite sites assumption: a character changes state once

$n$ mutations

$$F = \begin{pmatrix} 0.8 & 0.8 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.6 & 0.6 & 0.0 & 0.6 & 0.0 & 0.0 \\ 0.8 & 0.0 & 0.0 & 0.0 & 0.6 & 0.4 \end{pmatrix} m \text{ samples}$$

taxa

$$= \begin{pmatrix} 0.0 & 0.0 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.6 & 0.0 & 0.0 \\ 0.2 & 0.0 & 0.0 & 0.0 & 0.2 & 0.4 \end{pmatrix} \text{samples}$$

Usage Matrix $U$

mutations

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

Complete Two-State Perfect Phylogeny $B / T$

1-1

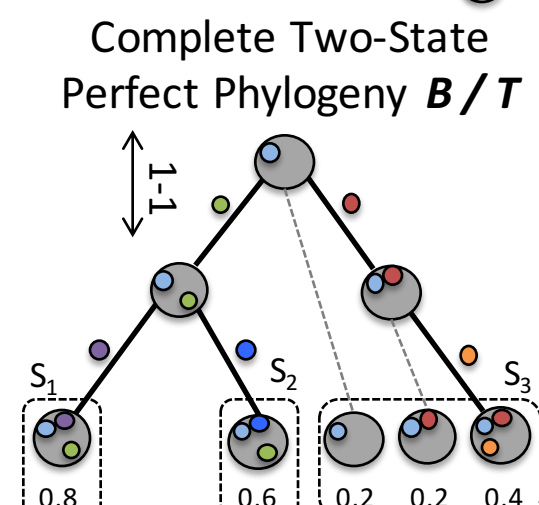$S_1$   $S_2$   $S_3$

0.8   0.6   0.2   0.2   0.4

$U = [u_{pj}]$ is a **usage matrix** iff

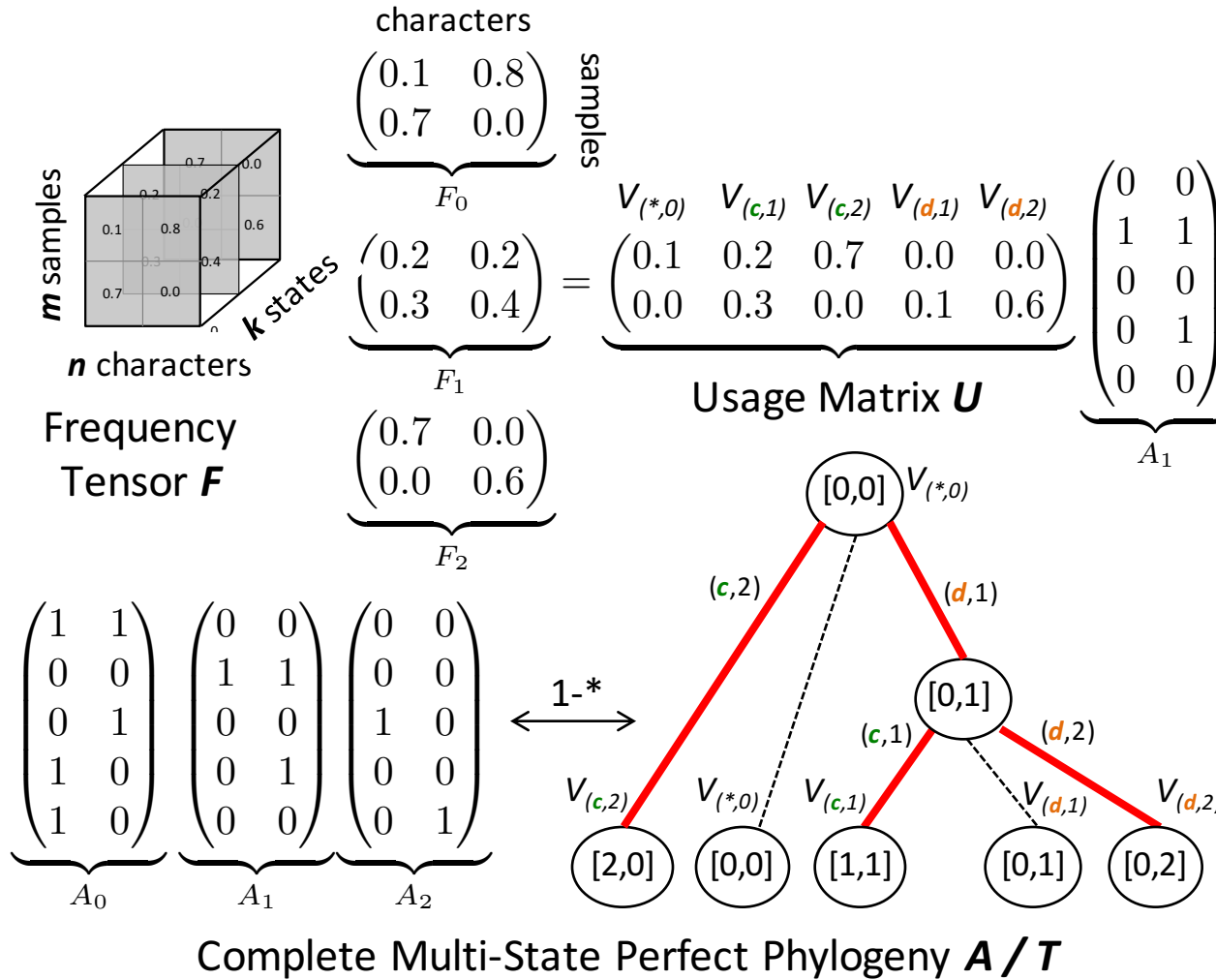$$u_{pj} \geq 0 \text{ and } \sum_j u_{pj} \leq 1$$

**VAF Factorization Problem (VAFFP):** [El-Kebir, Oesper et al., 2015]
Given $F$, find $U$ and $B$ such that $F = U B$

## Multi-State Perfect Phylogeny:

Infinite alleles assumption: a character changes to a state once

characters

$$F_0 = \begin{pmatrix} 0.1 & 0.8 \\ 0.7 & 0.0 \end{pmatrix} \text{samples}$$

$m$ samples

$n$ characters   $k$ states

Frequency Tensor $F$

$$F_1 = \begin{pmatrix} 0.2 & 0.2 \\ 0.3 & 0.4 \end{pmatrix}$$

$$F_2 = \begin{pmatrix} 0.7 & 0.0 \\ 0.0 & 0.6 \end{pmatrix}$$

$$\begin{array}{ccccc} V_{(*,0)} & V_{(c,1)} & V_{(c,2)} & V_{(d,1)} & V_{(d,2)} \end{array}$$

$$= \begin{pmatrix} 0.1 & 0.2 & 0.7 & 0.0 & 0.0 \\ 0.0 & 0.3 & 0.0 & 0.1 & 0.6 \end{pmatrix} \begin{pmatrix} 0 & 0 \\ 1 & 1 \\ 0 & 0 \\ 0 & 1 \\ 0 & 0 \end{pmatrix}$$

Usage Matrix $U$   $A_1$

$$A_0 = \begin{pmatrix} 1 & 1 \\ 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 0 \end{pmatrix} \quad A_1 = \begin{pmatrix} 0 & 0 \\ 1 & 1 \\ 0 & 0 \\ 0 & 1 \\ 0 & 0 \end{pmatrix} \quad A_2 = \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix}$$

1-*

[0,0] $V_{(*,0)}$

$(c,2)$   $(d,1)$

[0,1]

$(c,1)$   $(d,2)$

$V_{(c,2)}$   $V_{(*,0)}$   $V_{(c,1)}$   $V_{(d,1)}$   $V_{(d,2)}$

[2,0]   [0,0]   [1,1]   [0,1]   [0,2]

Complete Multi-State Perfect Phylogeny $A / T$

**Perfect Phylogeny Mixture Deconvolution Problem (PPMDP)**
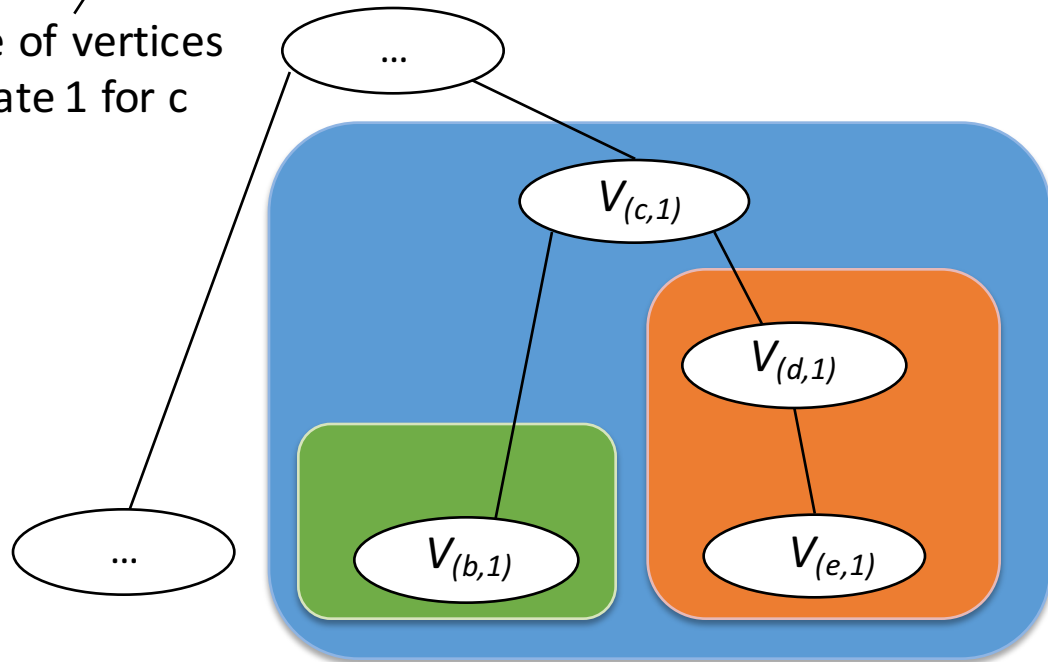[El-Kebir et al., 2016]: Given $F$, find $U$ and $A$ such that $F_i = U A_i$ for all states $i$

4

# Combinatorial Characterization

**Two-State Perfect Phylogeny:**

- A character changes state once
  - Once a mutation happens it persists
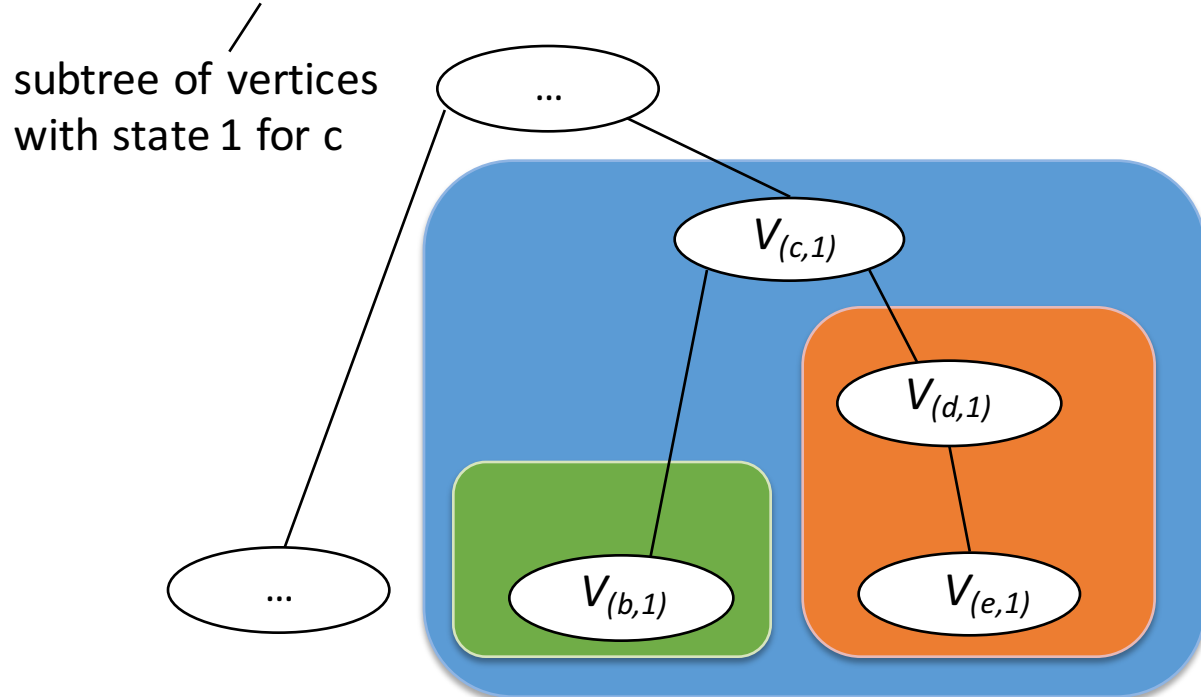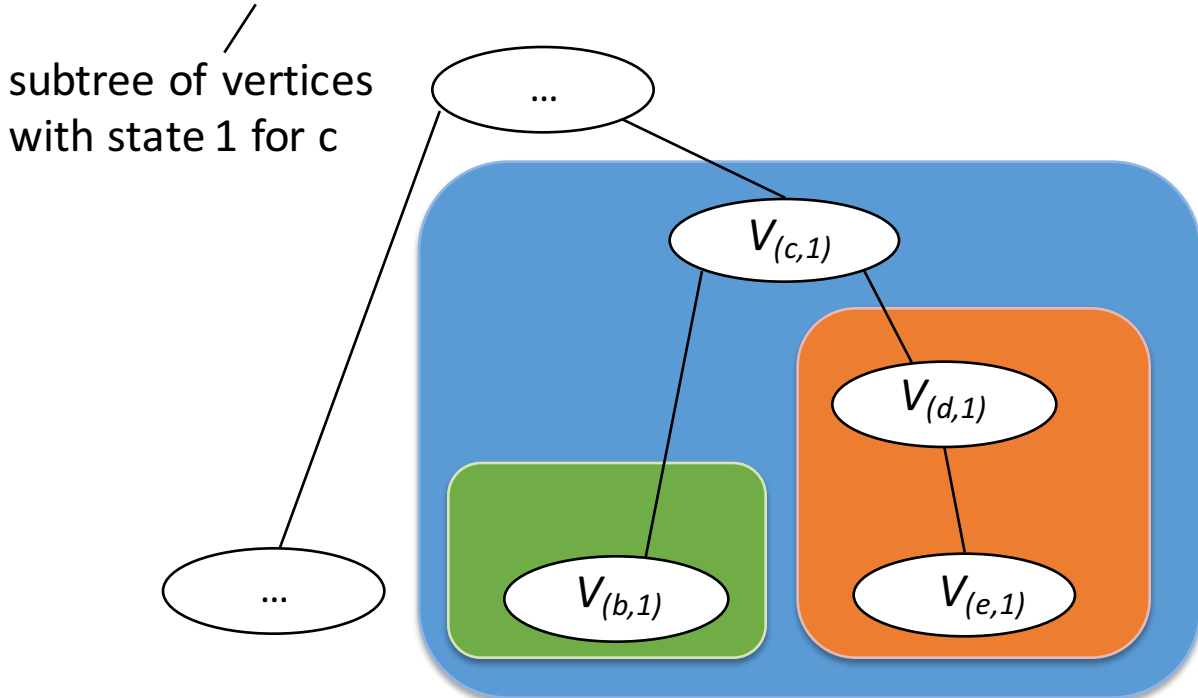- Thus $T_{(c,1)} = \bar{T}_{(c,1)}$ ————— subtree rooted at $V_{(c,1)}$

subtree of vertices
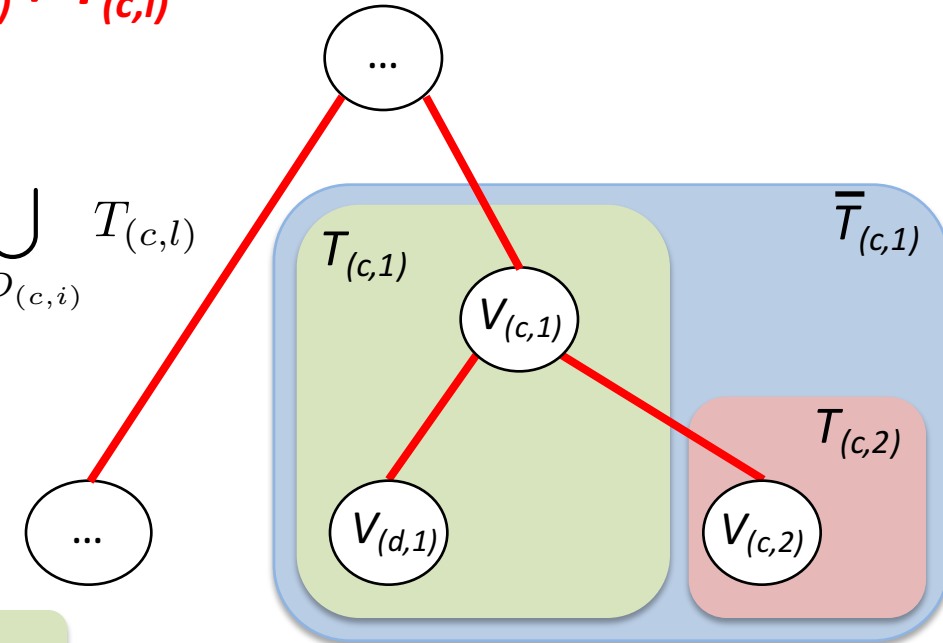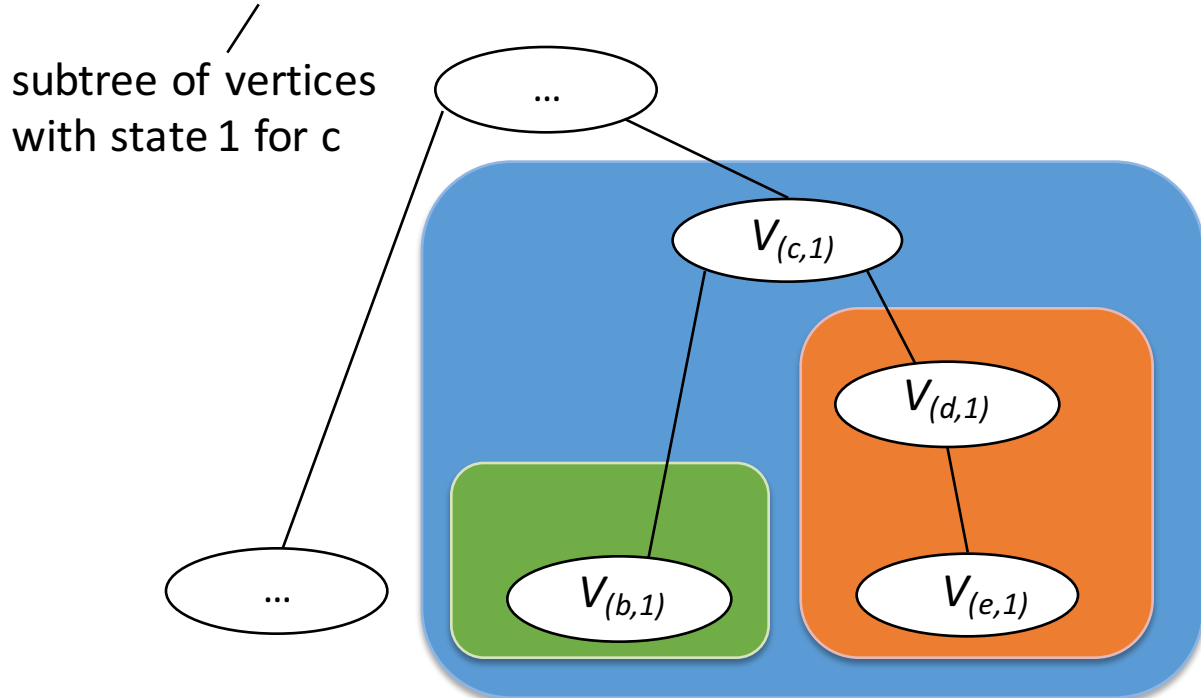with state 1 for c

# Combinatorial Characterization

**Two-State Perfect Phylogeny:**

- A character changes state once
  - Once a mutation happens it persists
- Thus $T_{(c,1)} = \bar{T}_{(c,1)}$ —— subtree rooted at $V_{(c,1)}$

subtree of vertices
with state 1 for c



**Sum Condition (SC)**

$$f_{p,(c,1)} \geq \sum_{(d,1) \in \delta(c,1)} f_{p,(d,1)}$$

# Combinatorial Characterization

**Two-State Perfect Phylogeny:**

- A character changes state once
  - Once a mutation happens it persists
- Thus $T_{(c,1)} = \overline{T}_{(c,1)}$ —— subtree rooted at $V_{(c,1)}$

subtree of vertices with state 1 for c
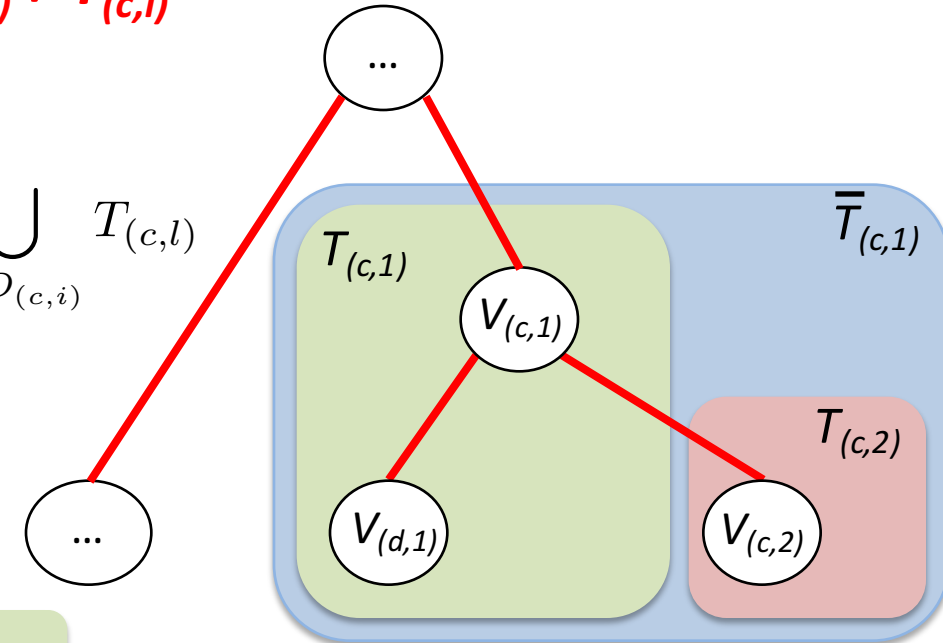


**Sum Condition (SC)**

$$f_{p,(c,1)} \geq \sum_{(d,1) \in \delta(c,1)} f_{p,(d,1)}$$

**Multi-State Perfect Phylogeny:**

- A character changes <u>to a state</u> once
- Thus, $T_{(c,i)} \neq \overline{T}_{(c,i)}$

- Instead:
$$\overline{T}_{(c,i)} = \bigcup_{l \in D_{(c,i)}} T_{(c,l)}$$



$$\overline{T}_{(c,1)} \neq T_{(c,1)}$$

$$\overline{T}_{(c,1)} = T_{(c,1)} \cup T_{(c,2)}$$

**Descendant set**
$D_{(c,1)} = \{1,2\}$

# Combinatorial Characterization

**Two-State Perfect Phylogeny:**

- A character changes state once
  - Once a mutation happens it persists
- Thus $T_{(c,1)} = \overline{T}_{(c,1)}$ —— subtree rooted at $V_{(c,1)}$

subtree of vertices with state 1 for c



**Multi-State Perfect Phylogeny:**

- A character changes <u>to a state</u> once
- Thus, $T_{(c,i)} \neq \overline{T}_{(c,i)}$

- Instead:
$$\overline{T}_{(c,i)} = \bigcup_{l \in D_{(c,i)}} T_{(c,l)}$$



$\overline{T}_{(c,1)} \neq T_{(c,1)}$

$\overline{T}_{(c,1)} = T_{(c,1)} \cup T_{(c,2)}$

**Descendant set**
$D_{(c,1)} = \{1,2\}$

**Sum Condition (SC)**
$$f_{p,(c,1)} \geq \sum_{(d,1) \in \delta(c,1)} f_{p,(d,1)}$$

**Multi-State Sum Condition (MSSC)** [El-Kebir et al., 2016]
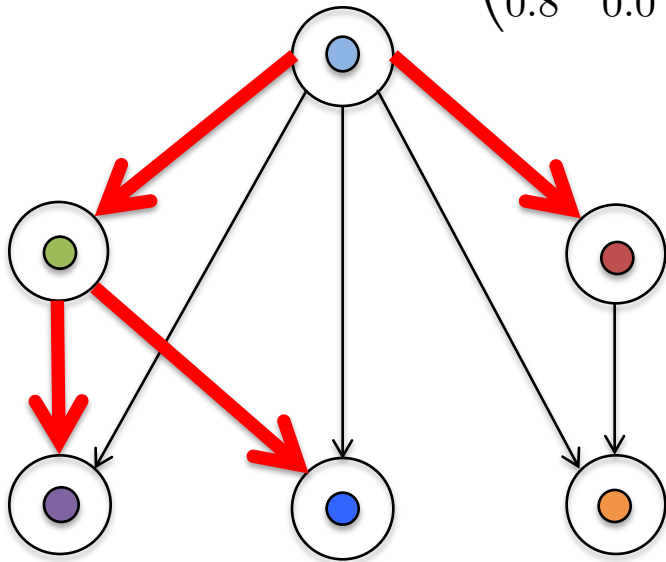
cumulative frequency $\longrightarrow$ $f_p^+(D_{(c,i)}) \geq \sum_{(d,j) \in \delta(c,i)} f_p^+(D_{(d,j)})$

5

# Spanning Trees in Ancestry Graph

**Two-State Perfect Phylogeny:**

mutations

$$\begin{pmatrix} 0.8 & 0.8 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.6 & 0.6 & 0.0 & 0.6 & 0.0 & 0.0 \\ 0.8 & 0.0 & 0.0 & 0.0 & 0.6 & 0.4 \end{pmatrix}$$ samples
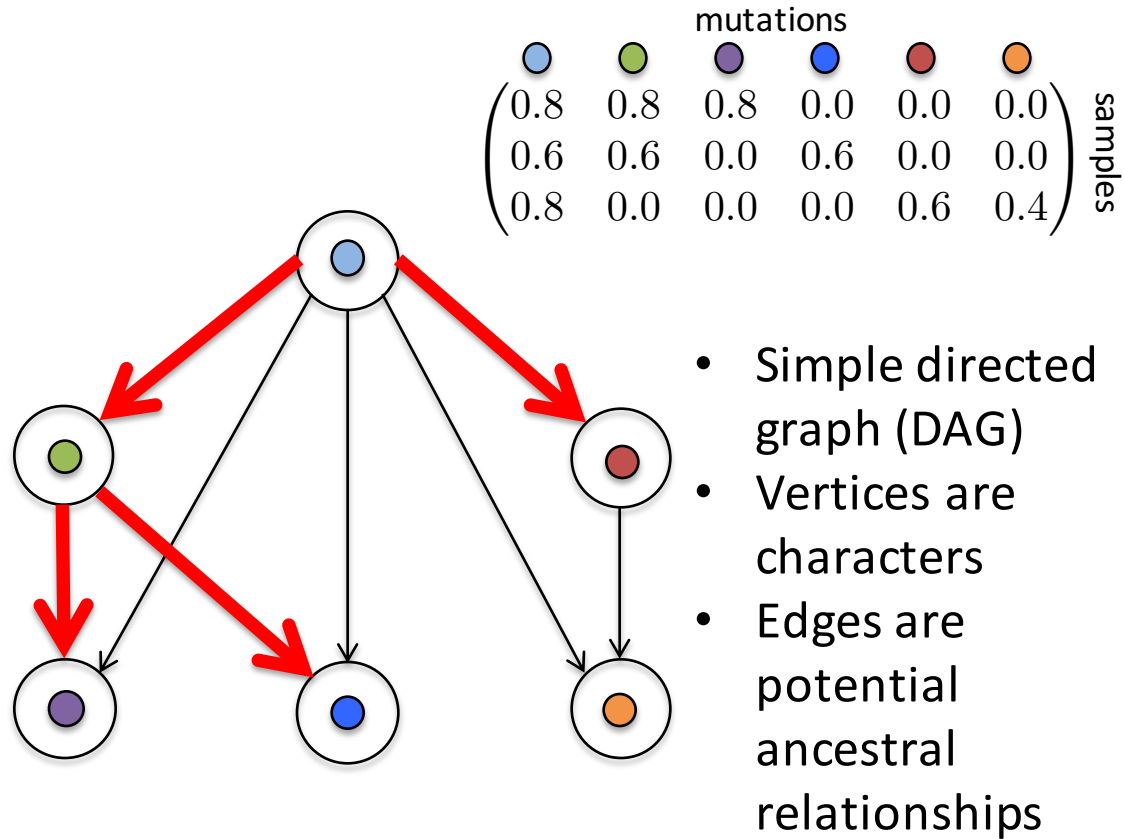
- Simple directed graph (DAG)
- Vertices are characters
- Edges are potential ancestral relationships

**Theorem 1** [El-Kebir, Oesper et al., 2015; Popic et al., 2015]
Solutions are spanning trees that satisfy (SC)

**Theorem 2** [El-Kebir, Oesper et al., 2015]
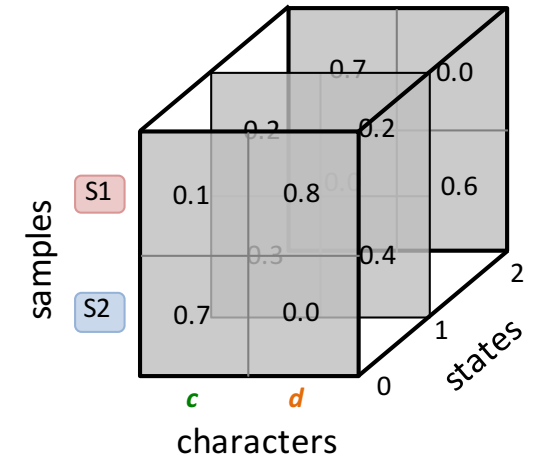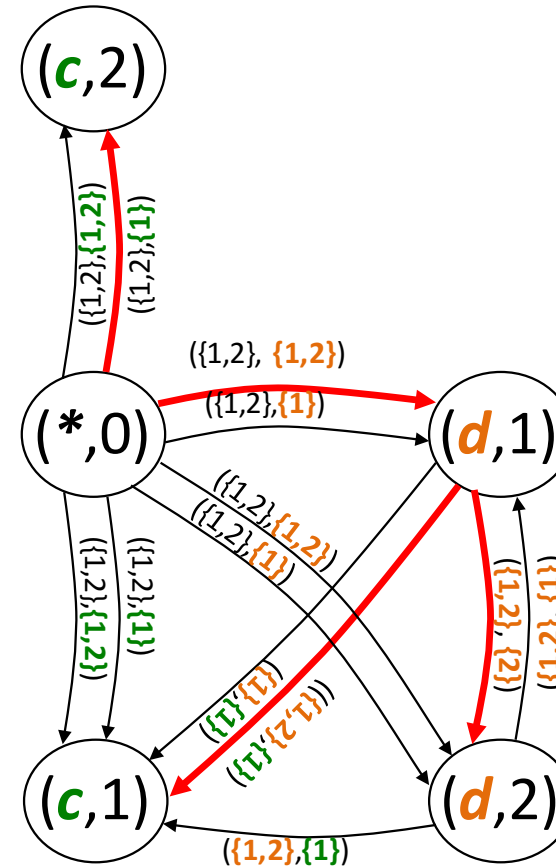VAFFP is NP-complete for $m = O(n)$

# Spanning Trees in Ancestry Graph

**Two-State Perfect Phylogeny:**



- Simple directed graph (DAG)
- Vertices are characters
- Edges are potential ancestral relationships

**Multi-State Perfect Phylogeny:**



- Directed multi-graph
- Vertices are character-state pairs
- Edges are labeled by valid descendant set pairs

**Theorem 1** [El-Kebir, Oesper et al., 2015; Popic et al., 2015]
Solutions are spanning trees that satisfy (SC)

**Theorem 2** [El-Kebir, Oesper et al., 2015]
VAFFP is NP-complete for $m = O(n)$

**Theorem 1** [El-Kebir et al., 2016]
Solutions are *threaded* spanning trees satisfying (MSSC)

**Theorem 2** [El-Kebir et al., 2016]
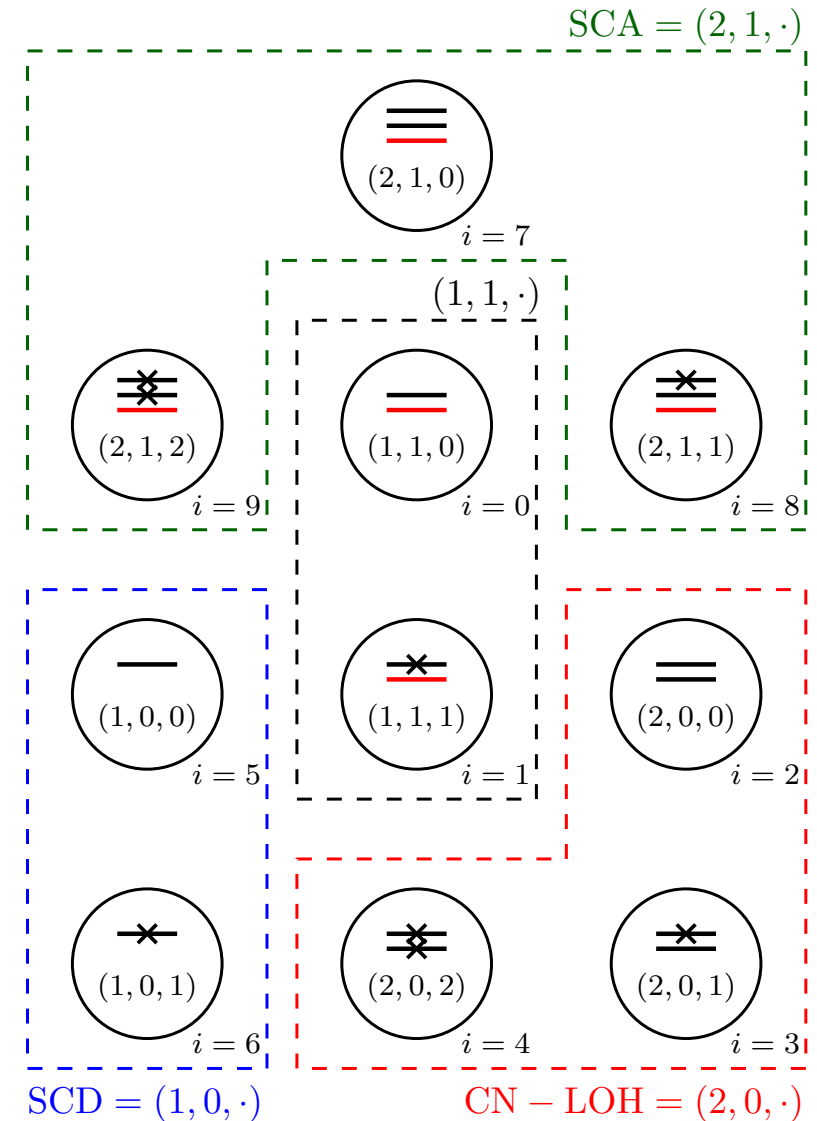PPMDP is NP-complete even for $m = 2$ and $k = 2$

# Application to Cancer Sequencing

**Input**

- Read-depth ratio
- B-allele frequencies
- Variant allele frequencies

**Model**

- Character is a genomic position (SNV)
- State is a triple (*x, y, z*) where
  - *x* is # maternal copies
  - *y* is # paternal copies
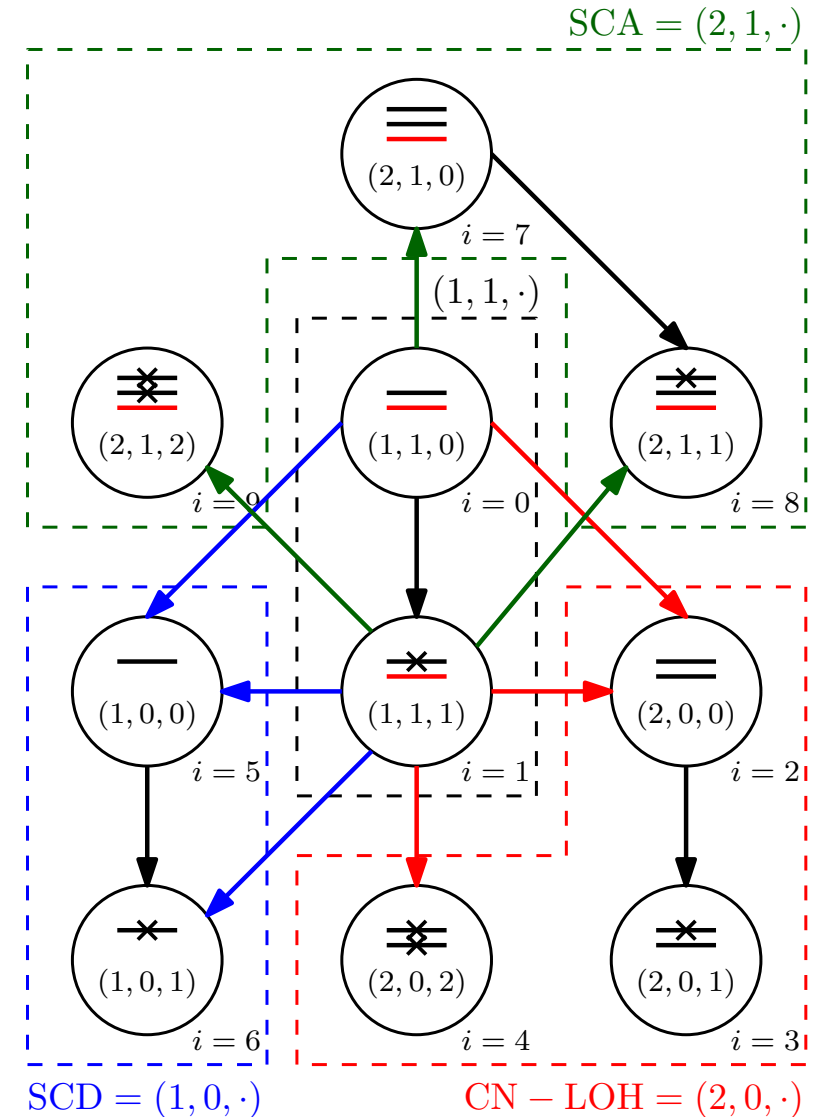  - *z* is # mutated copies
- Cladistic characters

# Application to Cancer Sequencing

**Input**

- Read-depth ratio
- B-allele frequencies
- Variant allele frequencies

**Model**

- Character is a genomic position (SNV)
- State is a triple (*x, y, z*) where
  - *x* is # maternal copies
  - *y* is # paternal copies
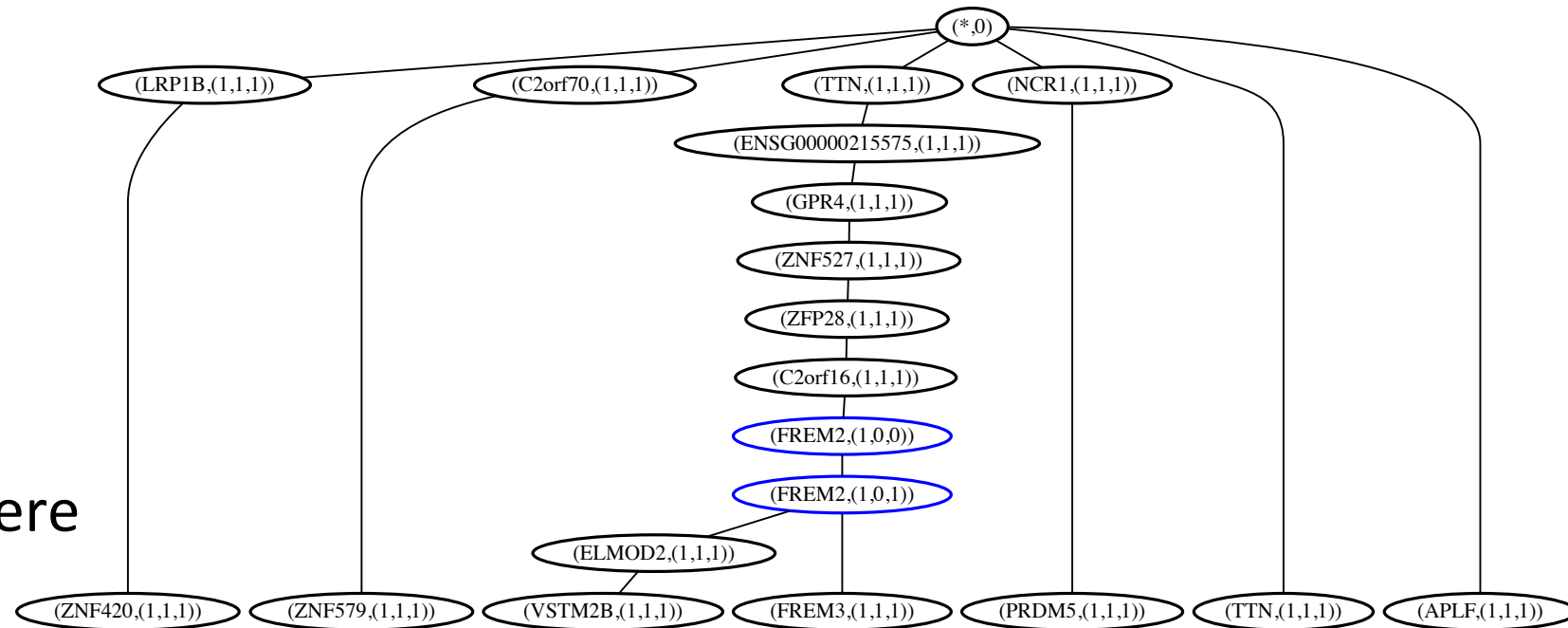  - *z* is # mutated copies
- Cladistic characters

# Application to Cancer Sequencing

**Input**

- Read-depth ratio
- B-allele frequencies
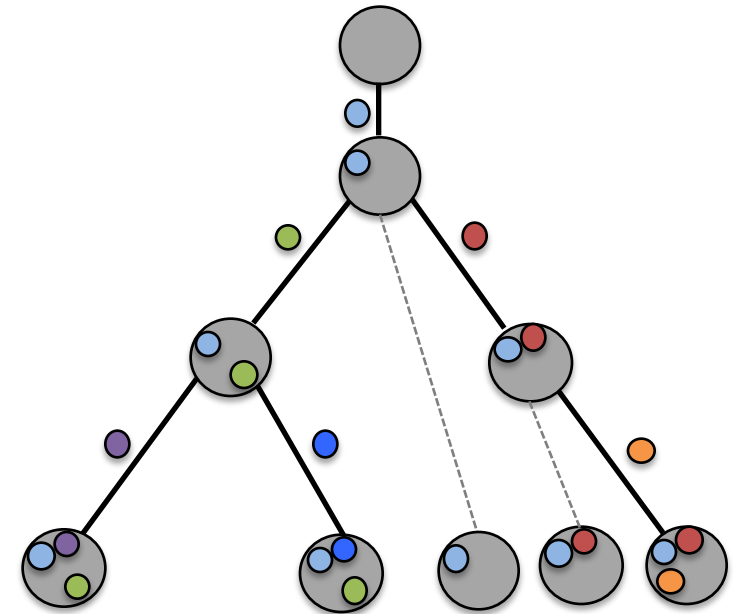- Variant allele frequencies

**Model**

- Character is a genomic position (SNV)
- State is a triple (*x, y, z*) where
  - *x* is # maternal copies
  - *y* is # paternal copies
  - *z* is # mutated copies
- Cladistic characters

# Conclusions

- Generalization of infinite sites model for SNVs is infinite alleles model for SNVs + CNAs

- Introduced Perfect Phylogeny Mixture Deconvolution Problem (PPMDP) for multi-state characters

- Combinatorial characterization of solutions

- PPMDP is NP-complete for $k = 2$ and $m = 2$

- Application to cancer sequencing
  - Metagenomics, somatic hypermutations, mtDNA, …

# Acknowledgements

**Research Group**

*Benjamin J. Raphael*
*Gryte Satas*
*Layla Oesper*
Dora Erdos
Matthew Reyna
Ashley Conard
Cyrus Cousins
Rebecca Elyanow
Hsin-Ta Wu

**Funding**

CCMB

**Preprint will be available soon on arXiv**