

Summarizing the Solution Space in Tumor Phylogeny Inference by Multiple Consensus Trees

Nuraini Aguse^{*}, Yuanyuan Qi^{*} and Mohammed El-Kebir University of Illinois at Urbana Champaign, Department of Computer Science

*Joint first authorship

Outline

- Introduction to tumor phylogenies
- Previous work
- Multiple Consensus Trees problem
- Methods
 - MILP
 - Heuristic
 - Bayesian Information Criterion



Introduction to tumor phylogenies

Cancer results from an evolutionary process



Tumor progression can be represented as a phylogenetic tree





Phylogenetic tree

Tumor progression can be represented as a phylogenetic tree





Phylogenetic tree

Tumor progression can be represented as a phylogenetic tree (or mutation tree)



1-1

Infinite Sites Assumption (ISA) states that a mutation is gained only once and never subsequently lost

Under ISA, a phylogenetic tree can be represented as a mutation tree

Phylogenetic tree

Mutation tree

Phylogenies enable us to understand and treat cancer



Phylogenies enable us to understand and treat cancer





- Bulk sequencing provides us frequencies of mutations as VAF
- Single-cell sequencing has a high false positive and false negative rate
- Both contribute to non-uniqueness of tumor phylogeny inference





- Bulk sequencing provides us frequencies of mutations as VAF
- Single-cell sequencing has a high false positive and false negative rate
- Both contribute to non-uniqueness of tumor phylogeny inference



Downstream analyses rely on accurate phylogeny inference



We need to summarize the solution space to

- Reduce inference errors
- Identify dependencies among mutations





We need to summarize the solution space to

- Reduce inference errors
- Identify dependencies among mutations

How do we best summarize the solution space?



Previous work

Solution space of lung cancer patient CRUK0037

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

Jamal-Hanjani et al. inferred 17 trees for patient CRUK0037



Solution space of lung cancer patient CRUK0037

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

Jamal-Hanjani et al. inferred 17 trees for patient CRUK0037



Two summary methods in previous work: Parent-child graph & single consensus tree





| | $v_1 \rightarrow v_{10}$ | $v_4 \rightarrow v_{10}$ |
|-----------------------|--------------------------|--------------------------|
| $v_1 \rightarrow v_7$ | 2 | 0 |
| $v_4 \rightarrow v_7$ | 2 | 5 |



| | $v_1 \rightarrow v_{10}$ | $v_4 \rightarrow v_{10}$ |
|-----------------------|--------------------------|--------------------------|
| $v_1 \rightarrow v_7$ | 2 | |
| $v_4 \rightarrow v_7$ | 2 | 5 |



| | $v_1 \rightarrow v_{10}$ | $v_4 \rightarrow v_{10}$ |
|-----------------------|--------------------------|--------------------------|
| $v_1 \rightarrow v_7$ | 2 | |
| $v_4 \rightarrow v_7$ | 2 | 5 |

Patterns of <u>mutual exclusivity</u> are not captured in parent-child graph

Summarizing the solution space using a single consensus tree



| | $v_1 \rightarrow v_{10}$ | $v_4 \rightarrow v_{10}$ |
|-----------------------|--------------------------|--------------------------|
| $v_1 \rightarrow v_7$ | 2 | 0 |
| $v_4 \rightarrow v_7$ | 2 | 5 |

Summarizing the solution space using a single consensus tree



| | $v_1 \rightarrow v_{10}$ | $v_4 \rightarrow v_{10}$ |
|-----------------------|--------------------------|--------------------------|
| $v_1 \rightarrow v_7$ | 2 | 0 |
| $v_4 \rightarrow v_7$ | 2 | 5 |

Summarizing the solution space using a single consensus tree



| | $v_1 \rightarrow v_{10}$ | $v_4 \rightarrow v_{10}$ |
|-----------------------|--------------------------|--------------------------|
| $v_1 \rightarrow v_7$ | 2 | 0 |
| $v_4 \rightarrow v_7$ | 2 | 5 |

Single consensus tree results in inaccurate summary of diverse solution spaces

Multiple Consensus Trees problem

Multiple Consensus Trees problem

Simultaneous clustering and consensus tree inference



Multiple Consensus Trees (MCT): [ISMB/ECCB 2019] Given trees $\mathcal{T} = \{T_1, ..., T_n\}$ and k > 0, find surjective clustering $\sigma : [n] \rightarrow [k]$ and consensus trees $\mathcal{R} = \{R_1, ..., R_k\}$ s.t. $\sum_{i=1}^n d(T_i, R_{\sigma(i)})$ is minimum⁵

Single Consensus Trees (SCT): [Govek et al., ACM-BCB 2018] Given $\mathcal{T} = \{T_1, ..., T_n\}$, find consensus tree R s.t. $\sum_{i=1}^n d(T_i, R)$ is minimum



Solution Space ${\mathcal T}$





Solution Space ${\mathcal T}$

Parent-Child distance

Single Consensus Trees (SCT): [Govek et al., ACM-BCB 2018] Given $\mathcal{T} = \{T_1, ..., T_n\}$, find consensus tree R s.t. $\sum_{i=1}^n d(T_i, R)$ is minimum

> **Theorem:** [Govek et al., ACM-BCB 2018] Max weight spanning arborescences of parent-child graph G_T are solutions to SCT



Solution Space ${\mathcal T}$

Single Consensus Trees (SCT): [Govek et al., ACM-BCB 2018] Given $\mathcal{T} = \{T_1, ..., T_n\}$, find consensus tree R s.t. $\sum_{i=1}^n d(T_i, R)$ is minimum

> **Theorem:** [Govek et al., ACM-BCB 2018] Max weight spanning arborescences of parent-child graph G_T are solutions to SCT





Single Consensus Trees (SCT): [Govek et al., ACM-BCB 2018] Given $\mathcal{T} = \{T_1, ..., T_n\}$, find consensus tree R s.t. $\sum_{i=1}^n d(T_i, R)$ is minimum

> **Theorem:** [Govek et al., ACM-BCB 2018] Max weight spanning arborescences of parent-child graph G_T are solutions to SCT

Multiple Consensus Trees (MCT): [Aguse et al., ISMB 2019] Given $\mathcal{T} = \{T_1, \dots, T_n\}$ and k > 0, find surjective clustering $\sigma : [n] \rightarrow [k]$ and consensus trees $\mathcal{R} = \{R_1, \dots, R_k\}$ s.t. $\sum_{i=1}^n d(T_i, R_{\sigma(i)})$ is minimum



Solution Space ${\mathcal T}$

Single Consensus Trees (SCT): [Govek et al., ACM-BCB 2018] Given $\mathcal{T} = \{T_1, ..., T_n\}$, find consensus tree R s.t. $\sum_{i=1}^n d(T_i, R)$ is minimum

> **Theorem:** [Govek et al., ACM-BCB 2018] Max weight spanning arborescences of parent-child graph G_T are solutions to SCT

Multiple Consensus Trees (MCT): [Aguse et al., ISMB 2019] Given $\mathcal{T} = \{T_1, \dots, T_n\}$ and k > 0, find surjective clustering $\sigma : [n] \rightarrow [k]$ and consensus trees $\mathcal{R} = \{R_1, \dots, R_k\}$ s.t. $\sum_{i=1}^n d(T_i, R_{\sigma(i)})$ is minimum

Proposition: [Aguse et al., ISMB 2019] Given fixed clustering $\sigma : [n] \rightarrow [k]$, MCT decomposes into k independent SCT instances



Solution Space ${\mathcal T}$

Single Consensus Trees (SCT): [Govek et al., ACM-BCB 2018] Given $\mathcal{T} = \{T_1, ..., T_n\}$, find consensus tree R s.t. $\sum_{i=1}^n d(T_i, R)$ is minimum

> **Theorem:** [Govek et al., ACM-BCB 2018] Max weight spanning arborescences of parent-child graph G_T are solutions to SCT

Multiple Consensus Trees (MCT): [Aguse et al., ISMB 2019] Given $\mathcal{T} = \{T_1, ..., T_n\}$ and k > 0, find surjective clustering $\sigma : [n] \rightarrow [k]$ and consensus trees $\mathcal{R} = \{R_1, ..., R_k\}$ s.t. $\sum_{i=1}^n d(T_i, R_{\sigma(i)})$ is minimum where $R_{\sigma(i)}$ is max weight spanning arborescence of $G_{\mathcal{T}_{\sigma(i)}}$

Proposition: [Aguse et al., ISMB 2019] Given fixed clustering $\sigma : [n] \rightarrow [k]$, MCT decomposes into k independent SCT instances



Solution Space ${\mathcal T}$

Single Consensus Trees (SCT): [Govek et al., ACM-BCB 2018] Given $\mathcal{T} = \{T_1, ..., T_n\}$, find consensus tree R s.t. $\sum_{i=1}^n d(T_i, R)$ is minimum

> **Theorem:** [Govek et al., ACM-BCB 2018] Max weight spanning arborescences of parent-child graph G_T are solutions to SCT

Multiple Consensus Trees (MCT): [Aguse et al., ISMB 2019] Given $\mathcal{T} = \{T_1, ..., T_n\}$ and k > 0, find surjective clustering $\sigma : [n] \rightarrow [k]$ and consensus trees $\mathcal{R} = \{R_{\pm}, ..., R_k\}$ s.t. $\sum_{i=1}^n d(T_i, R_{\sigma(i)})$ is minimum where $R_{\sigma(i)}$ is max weight spanning arborescence of $G_{\mathcal{T}_{\sigma(i)}}$

Proposition: [Aguse et al., ISMB 2019] Given fixed clustering $\sigma : [n] \rightarrow [k]$, MCT decomposes into k independent SCT instances



Single Consensus Trees (SCT): [Govek et al., ACM-BCB 2018] Given $\mathcal{T} = \{T_1, ..., T_n\}$, find consensus tree R s.t. $\sum_{i=1}^n d(T_i, R)$ is minimum

> **Theorem:** [Govek et al., ACM-BCB 2018] Max weight spanning arborescences of parent-child graph G_T are solutions to SCT

Multiple Consensus Trees (MCT): [Aguse et al., ISMB 2019] Given $\mathcal{T} = \{T_1, ..., T_n\}$ and k > 0, find surjective clustering $\sigma : [n] \rightarrow [k]$ and consensus trees $\mathcal{R} = \{R_{\pm}, ..., R_k\}$ s.t. $\sum_{i=1}^n d(T_i, R_{\sigma(i)})$ is minimum where $R_{\sigma(i)}$ is max weight spanning arborescence of $G_{\mathcal{T}_{\sigma(i)}}$

Proposition: [Aguse et al., ISMB 2019] Given fixed clustering $\sigma : [n] \rightarrow [k]$, MCT decomposes into k independent SCT instances



Methods and results

Mixed Integer Linear Program

Theorem: MCT is NP-hard for general k (by reduction from CLIQUE).

$$\begin{split} \min n(m-1) &- \sum_{i=1}^{n} \sum_{s=1}^{k} \sum_{p=1}^{m} \sum_{q=1}^{m} w_{i,s,p,q} \\ \text{s.t.} \quad \sum_{s=1}^{k} x_{i,s} &= 1 & \forall i \in [n] \\ &\sum_{i=1}^{n} x_{i,s} \geq 1 & \forall s \in [k] \\ &\sum_{p=1}^{m} z_{s,p} = 1 & \forall s \in [k] \\ &\sum_{q=1}^{m} y_{s,p,q} = 1 - z_{s,p} & \forall s \in [k], p \in [m] \\ &y_{s,p,q} \leq b_{p,q} & \forall s \in [k], p \in [m] \\ &\sum_{(p,q) \in \delta^{-}(U)} y_{s,p,q} + \sum_{p \in U} z_{s,p} \geq 1 & \forall s \in [k], U \subseteq [m] \\ &w_{i,s,p,q} \leq a_{i,p,q} & \forall i \in [n], s \in [k], p, q \in [m] \\ &w_{i,s,p,q} \leq x_{i,s} & \forall i \in [n], s \in [k], p, q \in [m] \\ &w_{i,s,p,q} \leq y_{s,p,q} & \forall i \in [n], s \in [k], p, q \in [m] \\ &w_{i,s,p,q} \geq 0 & \forall i \in [n], s \in [k], p, q \in [m] \\ &w_{i,s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ &y_{s,p,q} \leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall$$

s.t.

$$\forall s \in [k], p, q \in [m]$$

$$\sum_{i=1}^{n} x_{i,s} \ge \sum_{i=1}^{n} x_{i,s+1} + 1 \qquad \forall s \in [k-1]$$
$$x_{i,s} \in \{0,1\} \qquad \forall i \in [n], s \in [k]$$

 $y_{s,p,q} \ge \sum_{i=1}^{n} a_{i,p,q} x_{i,s} - \sum_{i=1}^{n} x_{i,s} + 1$

 $y_{s,p,q} \ge 0$ $z_{s,p} \ge 0$

$$\forall s \in [k], p, q \in [m]$$

$$\forall s \in [k], p \in [m]$$

Mixed Integer Linear Program

Theorem: MCT is NP-hard for general k (by reduction from CLIQUE).

 $egin{aligned} x_{i,s} \in \{0,1\} & ext{Tree } T_i ext{ is assigned to cluster } s \ y_{s,p,q} \geq 0 & ext{Edge } (p,q) ext{ is present in consensus tree } R_s \ z_{s,p} \geq 0 & ext{Vertex } p ext{ is root of consensus tree } R_s \end{aligned}$

$$\begin{aligned} \min n(m-1) &- \sum_{i=1}^{n} \sum_{s=1}^{k} \sum_{p=1}^{m} \sum_{q=1}^{m} w_{i,s,p,q} \\ \text{s.t.} \quad \sum_{s=1}^{k} x_{i,s} &= 1 & \forall i \in [n] \\ \sum_{i=1}^{n} x_{i,s} &\geq 1 & \forall s \in [k] \\ \sum_{i=1}^{m} x_{i,s} &\geq 1 & \forall s \in [k] \\ \sum_{q=1}^{m} z_{s,p} &= 1 - z_{s,p} & \forall s \in [k], p \in [m] \\ y_{s,p,q} &\leq b_{p,q} & \forall s \in [k], p, q \in [m] \\ \sum_{(p,q) \in \delta^{-}(U)} y_{s,p,q} + \sum_{p \in U} z_{s,p} &\geq 1 & \forall s \in [k], U \subseteq [m] \\ w_{i,s,p,q} &\leq a_{i,p,q} & \forall i \in [n], s \in [k], p, q \in [m] \\ w_{i,s,p,q} &\leq x_{i,s} & \forall i \in [n], s \in [k], p, q \in [m] \\ w_{i,s,p,q} &\leq x_{i,s} & \forall i \in [n], s \in [k], p, q \in [m] \\ w_{i,s,p,q} &\leq y_{s,p,q} & \forall i \in [n], s \in [k], p, q \in [m] \\ w_{i,s,p,q} &\geq 0 & \forall i \in [n], s \in [k], p, q \in [m] \\ y_{s,p,q} &\leq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} & \forall s \in [k], p, q \in [m] \\ y_{s,p,q} &\geq \sum_{i=1}^{n} a_{i,p,q} x_{i,s} - \sum_{i=1}^{n} x_{i,s} + 1 & \forall s \in [k], p, q \in [m] \\ \sum_{i=1}^{n} x_{i,s} &\geq \sum_{i=1}^{n} x_{i,s+1} + 1 & \forall s \in [k], p, q \in [m] \\ x_{i,s} &\in \{0,1\} & \forall i \in [n], s \in [k], p, q \in [m] \\ y_{s,p,q} &\geq 0 & \forall s \in [k], p, q \in [m] \end{aligned}$$

 $\forall s \in [k], p \in [m]$

MILP does not scale well with k and n

50

small medium large 40 40 30 30 11 40 11 40 11 40 11 40 11 10 40 11 10 40 11 10 0 2 3 1 4 5 #clusters k

We simulated **small**, **medium** and **large** instances of phylogeny inference solution spaces

Coordinate Ascent heuristic



Coordinate Ascent heuristic finds optimal solutions efficiently

| | #clusters k | MILP (1 h) | CA (100 r.) |
|------|-------------|------------|-------------|
| () | 2 | 16 | 16 |
| (1 | 3 | 16 | 16 |
| lla | 4 | 16 | 16 |
| SIT | 5 | 16 | 16 |
| (5) | 2 | 15 | 15 |
| n (] | 3 | 13 | 13 |
| iun | 4 | 12 | 12 |
| led | 5 | 10 | 10 |
| 4)n | 2 | 3 | 3 |
| (1 | 3 | 0 | 0 |
| rge | 4 | 0 | 0 |
| laı | 5 | 0 | 0 |
| | | | - |

Small, medium, and large simulated instances

Number of instances solved by MILP to provable optimality

Coordinate Ascent heuristic finds optimal solutions efficiently

| | #clusters k | MILP (1 h) | CA (100 r.) |
|---------|-------------|------------|-------------|
| (9 | 2 | 16 | 16 |
| (1) | 3 | 16 | 16 |
| lla | 4 | 16 | 16 |
| sn | 5 | 16 | 16 |
| [5] | 2 | 15 | 15 |
| u (] | 3 | 13 | 13 |
| iun | 4 | 12 | 12 |
| led | 5 | 10 | 10 |
| 4) n | 2 | 3 | 3 |
| (1^7) | 3 | 0 | 0 |
| rge | 4 | 0 | 0 |
| laı | 5 | 0 | 0 |
| | | | |

Small, medium, and large simulated instances

Number of instances where heuristic returned MILP's optimal solution

Coordinate Ascent heuristic finds optimal solutions efficiently

| | #clusters k | MILP (1 h) | CA (100 r.) |
|------|-------------|------------|-------------|
| 6) | 2 | 16 | 16 |
| (1 | 3 | 16 | 16 |
| lla | 4 | 16 | 16 |
| SIT | 5 | 16 | 16 |
| 15) | 2 | 15 | 15 |
| n (] | 3 | 13 | 13 |
| iun | 4 | 12 | 12 |
| ned | 5 | 10 | 10 |
| 4)n | 2 | 3 | 3 |
| (17 | 3 | 0 | 0 |
| rge | 4 | 0 | 0 |
| laı | 5 | 0 | 0 |

Small, medium, and large simulated instances

What about the number of clusters, k?

Bayesian Information Criterion determines the number of clusters for each solution space

Jamal-Hanjani et al. (2017). NEJM.

Jamal-Hanjani et al. inferred 8 trees for patient CRUK0013





Bayesian Information Criterion determines the number of clusters for each solution space

Jamal-Hanjani et al. (2017). NEJM.

Jamal-Hanjani et al. inferred 17 trees for patient CRUK0037





Multiple Consensus Trees capture patterns of mutual exclusivity and co-occurrence



Multiple Consensus Trees capture patterns of mutual exclusivity and co-occurrence



These edges tend to co-occur in the trees in the solution space





Conclusions

- Introduced the Multiple Consensus Trees (MCT) problem
- Showed hardness and presented a mixed integer linear program
- Presented an efficient heuristic and showed that it finds optimal solutions
- Model selection for the number of clusters

Future directions

- Relax infinite sites assumption
- Use alternative distance functions



CENTER FOR COMPUTATIONAL BIOTECHNOLOGY & GENOMIC MEDICINE

Acknowledgements

- El-Kebir group
 - Mohammed El-Kebir
 - Yuanyuan Qi
 - Juho Kim
 - Jiaqi Wu
- ISCB for Travel Fellowship Award
- UIUC Center for Computational Biotechnology and Genomic Medicine (grant: CSN 1624790)
- National Science Foundation (CCF-1850502)

