

to s then $r'[i] = s$. Note that in such a case at least one of $M0, M1$ contain both the states on i and $i, c(v)$ share at least 3 gametes in M_j . The proof can be split into two symmetric cases based on whether r is fixed on condition 1 or 2. One case is presented below:

Taxon $r[i]$ is fixed based on condition 1: In this case, all the taxa in $M0$ contain the same state s on i . Therefore, the taxa in $M1$ should contain both states on i . Hence i mutates in $T_{M_j}^*(c(v), 1)$. For the sake of contradiction, assume that $r'[i] \neq s$. If $i \notin \mu(e)$ then $p'[i] \neq s$. However, all the taxa in $M0$ contain state s . This implies that i mutates in $T_{M_j}^*(c(v), 0)$ as well. Therefore i and $c(v)$ share all four gametes on $T_{M_j}^*$. However i and $c(v)$ share at most 3 gametes in M_j - one in $M0$ and at most two in $M1$. This leads to a contradiction to Lemma 4.5. Once r is guessed correctly, p can be computed since it is identical to r in all characters except $c(v)$ and those that share two gametes with $c(v)$ in M_j . We make a note here that we are assuming that e does not mutate any character that does not share two gametes with $c(v)$ in M_j . This creates a small problem that although the length of the tree constructed is optimal, r and p could be degree-two Steiner vertices. If after constructing the optimum phylogenies for $M0$ and $M1$, we realize that this is the case, then we simply add the mutation adjacent to r and p to the edge (r, p) and return the resulting phylogeny where both r and p are not degree-two Steiner vertices.

The above implementation therefore requires only guessing states corresponding to the remaining unfixed characters of r . If a character i violates the first two conditions, then i mutates once in $T_{M_j}^*(i, 0)$ and once in $T_{M_j}^*(i, 1)$. If $r[i]$ has not been fixed, then we can associate a pair of mutations of the same character i with it. At the end of the current iteration M_j is replaced with $M0$ and $M1$ and each contains exactly one of the two associated mutations. Therefore if q' characters are unfixed then $\text{penalty}(M0) + \text{penalty}(M1) \leq \text{penalty}(M_j) - q'$. Since $\text{penalty}(I) \leq q$, throughout the execution of the algorithm there are q unfixed states. Therefore the probability of all the guesses being correct is 2^{-q} . ■

This completes our analysis for upper bounding the probability that the algorithm returns an optimum phylogeny. We now analyze the running time. We use the following lemma to show that we can efficiently construct optimum phylogenies at Step 3 in the pseudo-code:

Lemma 4.7: For a set of taxa M , if the number of non-isolated vertices of the associated conflict graph is t , then an optimum phylogeny T_M^* can be constructed in time $O(3^s 6^t + nm^2)$, where $s = \text{penalty}(M)$.