

Theorem S1. Let Σ^* denote the complete set of sequences that are concatenations of motifs from Σ . If an efficient motif set $\tilde{\Sigma} \subseteq \Sigma$ exists, which minimizes the sum of $\lambda \times \|\tilde{\Sigma}\|$ and total motifs replacement cost $\sum_{i=1}^p o_i * \sum_{j=1}^p x_{ij} * \delta_{ij}$ and satisfies the three requirements below, then there exists $\tilde{v}_j \in \tilde{\Sigma}^*$ such that $\sum_j \text{div}(v_j, \tilde{v}_j) \leq \sum_{i=1}^p o_i * \sum_{j=1}^p x_{ij} * \delta_{ij} < \Delta, v_j \in V$.

- All occurrence of $m_i \in \Sigma$ is replaced by one and only one $m_j \in \tilde{\Sigma}$ (possibly itself).
- Motif m_i cannot be replaced by motif m_j if $o_i \geq o_j$.
- The total motifs replacement cost $\sum_{i=1}^p o_i * \sum_{j=1}^p x_{ij} * \delta_{ij} < \Delta$.

Proof. Assume an efficient motif set $\tilde{\Sigma} \subseteq \Sigma$ exists. Each VNTR sequence v_j can be represented as a sequence of original motifs $m_j^1 \circ \dots \circ m_j^l$, each $m_j^i \subseteq \Sigma$. By substituting each m_j^i with its counterpart efficient motif \tilde{m}_j^i , we get $\tilde{v}_j = \tilde{m}_j^1 \circ \dots \circ \tilde{m}_j^l$. For each v_j and \tilde{v}_j , it is clear that $\text{div}(v_j, \tilde{v}_j) \leq \sum_i \text{div}(m_j^i, \tilde{m}_j^i)$, otherwise the combination of edit operations from individual motif alignments would infer a more parsimonious edit distance for v_j and \tilde{v}_j . Therefore, $\sum_{j=1}^p \text{div}(v_j, \tilde{v}_j) \leq \sum_i \text{div}(m_j^i, \tilde{m}_j^i) = \sum_{i=1}^p o_i * \sum_{j=1}^p x_{ij} * \delta_{ij} < \Delta$. \square