

Multiple Sequence Alignment

Chin Lung Lu

Computational Biology

Analyses and Applications of Sequences

Multiple sequence alignment

S_1 : RCTLEE	S_1 : R C T L E E
S_2 : RCLEE	S_2 : R C - L E E
S_3 : CTLEE	S_3 : - C T L E E
S_4 : CTEE	S_4 : - C T - E E

4 Sequences \Rightarrow A Multiple Sequence Alignment

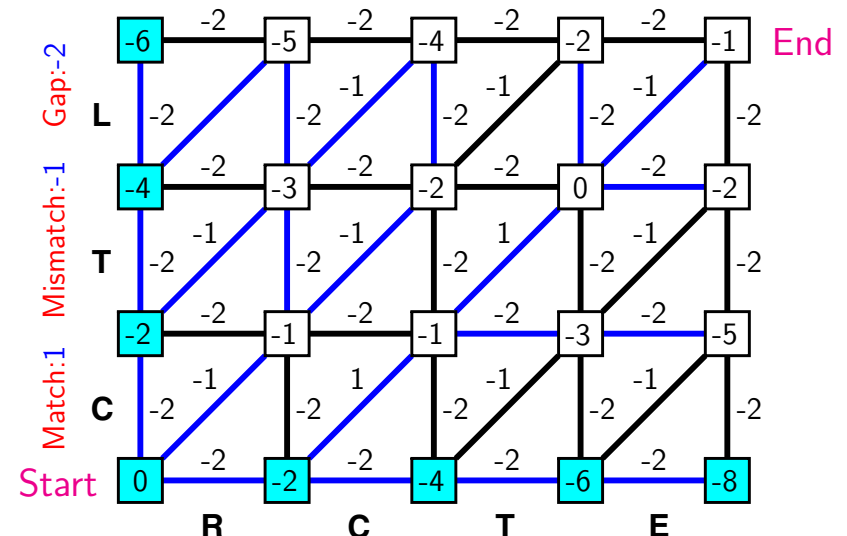
SP-score of an MSA

- **SP-score** (Sum-of-Pairs) of an MSA: the sum of the SP-scores of all columns
- **SP-score of a column**: the sum of pairwise scores of all pairs of characters

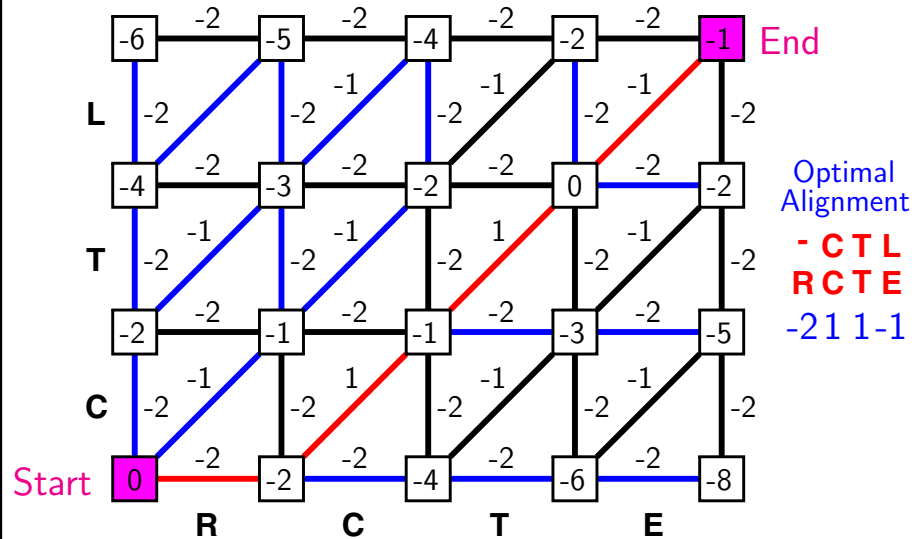
S_1 :	R	C	T	L	E	E
S_2 :	R	C	-	L	E	E
S_3 :	-	C	T	L	E	E
S_4 :	-	C	T	-	E	E

The SP-score of the first column is $score(R, R) + 4 \times score(R, -) + score(-, -)$, where $score(-, -) = 0$

Two sequence alignment



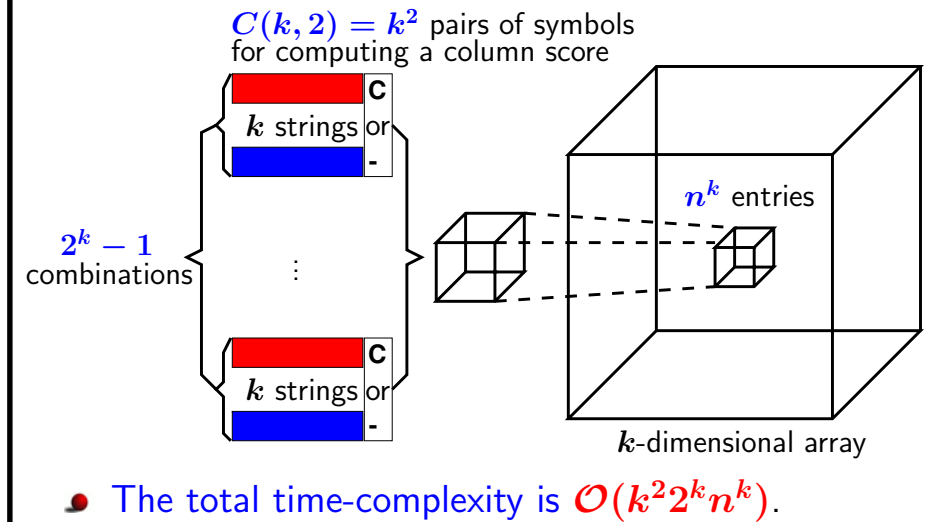
Two sequence alignment



By C.L. Lu

Multiple Sequence Alignment p.5

Sum-of-pairs MSA problem



By C.L. Lu

Multiple Sequence Alignment p.6

Sum-of-pairs MSA problem

- **NP-complete:** by Wang and Jiang [WJ94]; Bonizzoni and Vedova [BV01]
- **Branch and bound methods:** by Carrillo and Lipman [CL88] (improved by [LAK89, GKS95])
- **Approximate methods:** (k sequences)
 - $(2 - \frac{2}{k})$ -approximation, by Gusfield [Gus93]
 - $(2 - \frac{3}{k})$ -approximation, by Pevzner [Pev92]
 - $(2 - \frac{l}{k})$ -approximation, where $l < k$ by Bafna, Lawler and Pevzner [BLP97]
- **Heuristic methods:** Progressive algorithms

By C.L. Lu

Multiple Sequence Alignment p.7

Branch and bound strategy

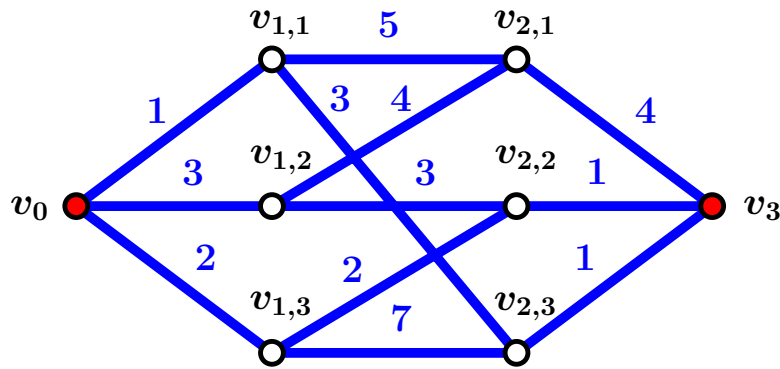
1. **Branch mechanism:** a way of generating the branches of the solution space (tree)
2. **Bound mechanism:** a way of bounding the branches to avoid the exhaustive search of solution space
 - Find an upper bound of an optimal solution (find a feasible solution using tree searching methods, like breadth-first, depth-first, and best-first searches)
 - Predict a lower bound for a branch
 - If the lower bound of a branch $>$ the upper bound, then the branch is terminated.

By C.L. Lu

Multiple Sequence Alignment p.8

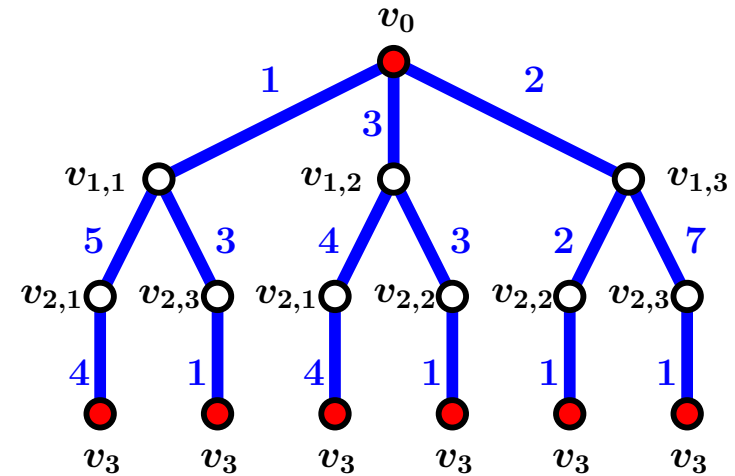
Shortest path problem ①

- Find a shortest path from v_0 to v_3

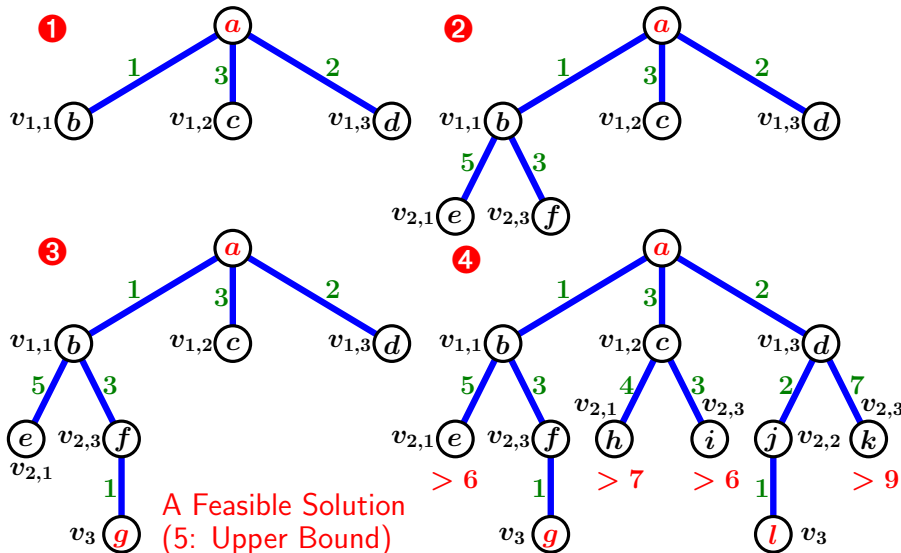


Shortest path problem ②

- A search tree of the problem (6 feasible solutions):



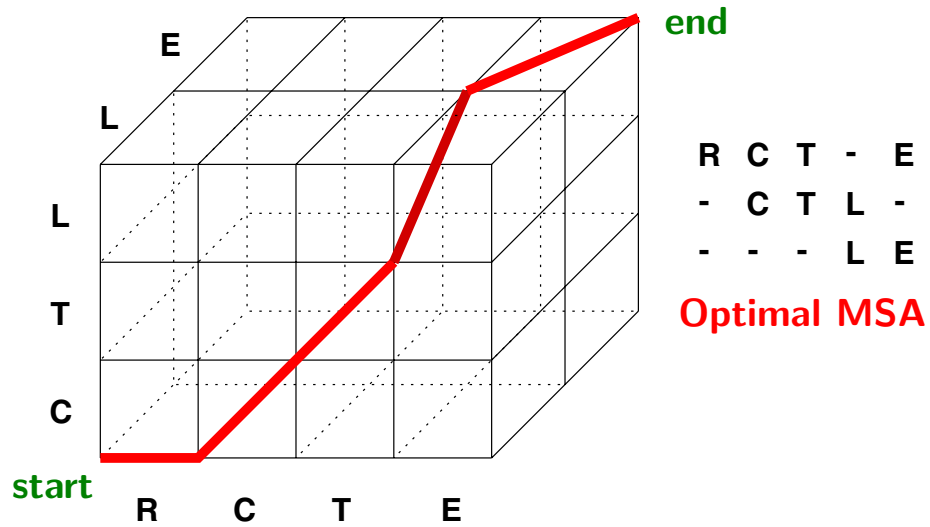
Shortest path problem ③



Sum-of-pairs MSA problem

- NP-complete: by Wang and Jiang [WJ94]; Bonizzoni and Vedova [BV01]
- Branch and bound methods: by Carrillo and Lipman [CL88] (improved by [LAK89, GKS95])
- Approximate methods: (k sequences)
 - $(2 - \frac{2}{k})$ -approximation, by Gusfield [Gus93]
 - $(2 - \frac{3}{k})$ -approximation, by Pevzner [Pev92]
 - $(2 - \frac{l}{k})$ -approximation, where $l < k$ by Bafna, Lawler and Pevzner [BLP97]
- Heuristic methods: Progressive algorithms

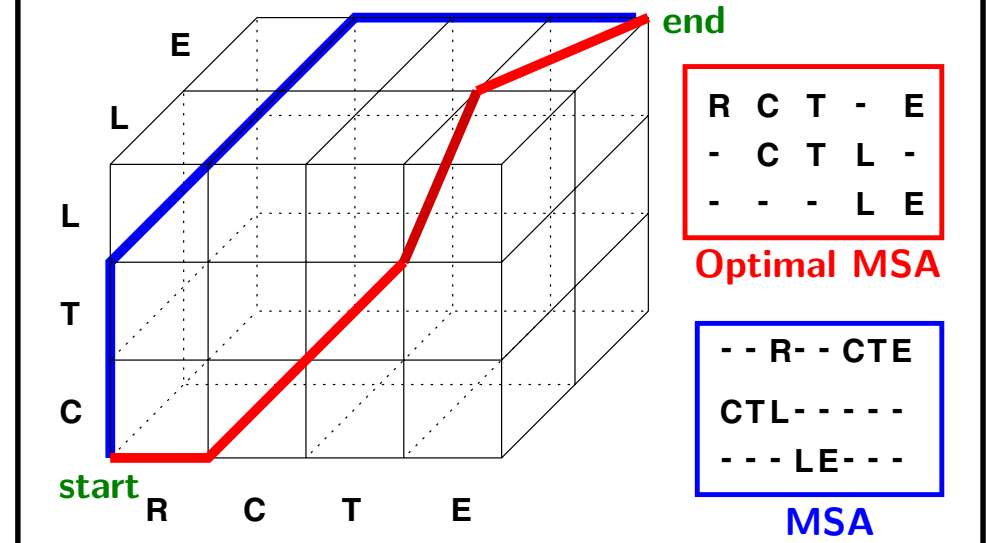
Dynamic programming approach



By C.L. Lu

Multiple Sequence Alignment p.13

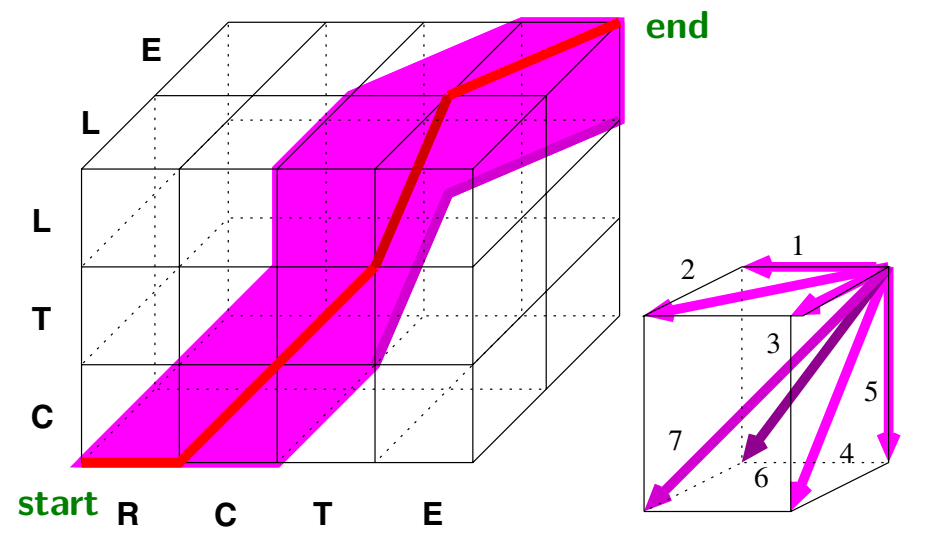
Dynamic programming approach



By C.L. Lu

Multiple Sequence Alignment p.14

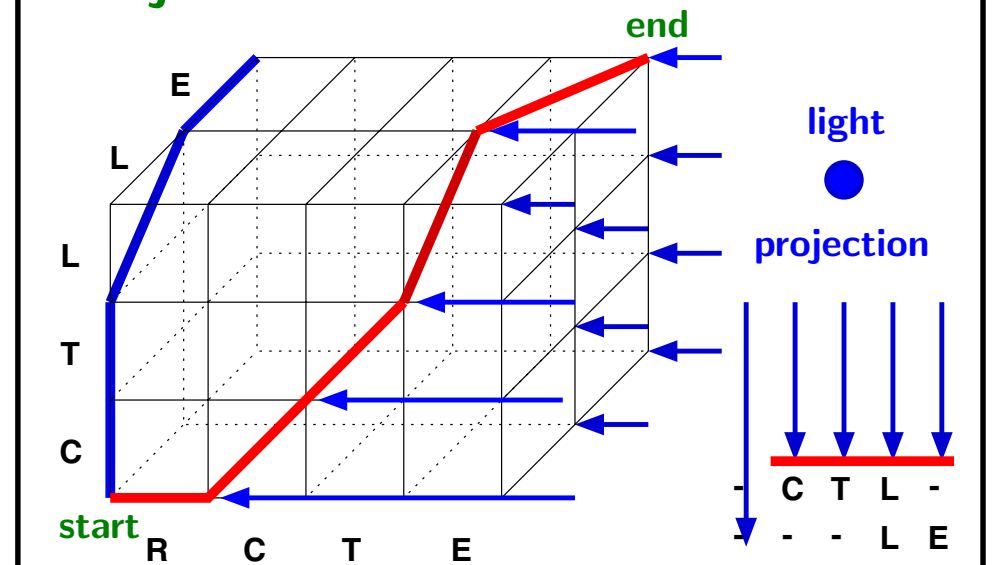
Polyhedron



By C.L. Lu

Multiple Sequence Alignment p.15

Projection of MSA



By C.L. Lu

Multiple Sequence Alignment p.16

Induced pairwise alignment

- Let \mathcal{A} be a multiple alignment of S_1, S_2, \dots, S_k .
- Induced pairwise alignment $\mathcal{A}_{i,j}$
(a projection of \mathcal{A} w.r.t. S_i and S_j):
a pairwise alignment of S_i and S_j obtained by copying rows i and j of \mathcal{A} and removing the columns with a gap in both rows
- $c(\mathcal{A}_{i,j}) \geq d(S_i, S_j)$
 - $c(\mathcal{A}_{i,j})$: the cost of $\mathcal{A}_{i,j}$
 - $d(S_i, S_j)$: the minimum distance of S_i and S_j
- The SP-cost of $\mathcal{A} = c(\mathcal{A}) = \sum_{i < j} c(\mathcal{A}_{i,j})$

Costing function: metric

- A costing function c is a **metric** if for any x, y and z , it satisfies ① $c(x, y) \geq 0$, where $c(x, y) = 0$ if $x = y$, ② **symmetry**: $c(x, y) = c(y, x)$ and ③ **triangle inequality**: $c(x, z) \leq c(x, y) + c(y, z)$.
- Given a set S of k sequences S_1, S_2, \dots, S_k and a scoring function c (**metric**):
 - \mathcal{A}^o : an optimal MSA of S with cost $c(\mathcal{A}^o)$
 - \mathcal{A}^h : a heuristic MSA of S with cost $c(\mathcal{A}^h)$
 - $\mathcal{A}_{i,j}^o$: the projection of \mathcal{A}^o w.r.t. S_i and S_j
 - $\mathcal{A}_{i,j}^h$: the projection of \mathcal{A}^h w.r.t. S_i and S_j

Carrillo-Lipman bound

- SP-cost(\mathcal{A}^o) = $c(\mathcal{A}^o) = \sum_{1 \leq i < j \leq k} c(\mathcal{A}_{i,j}^o)$
- SP-cost(\mathcal{A}^h) = $c(\mathcal{A}^h) = \sum_{1 \leq i < j \leq k} c(\mathcal{A}_{i,j}^h)$
- $c(\mathcal{A}^h) \geq c(\mathcal{A}^o)$
 - $\Rightarrow \sum_{1 \leq i < j \leq k} c(\mathcal{A}_{i,j}^h) \geq \sum_{1 \leq i < j \leq k} c(\mathcal{A}_{i,j}^o)$
 - $\Rightarrow \sum_{1 \leq i < j \leq k} (c(\mathcal{A}_{i,j}^h) - c(\mathcal{A}_{i,j}^o)) \geq 0$

Carrillo-Lipman bound

- $\because \sum_{1 \leq i < j \leq k} (c(\mathcal{A}_{i,j}^h) - c(\mathcal{A}_{i,j}^o)) \geq 0$
- \therefore For any arbitrary projection w.r.t. S_p and S_q ,
 - $\left(\sum_{\substack{1 \leq i < j \leq k \\ (i,j) \neq (p,q)}} (c(\mathcal{A}_{i,j}^h) - c(\mathcal{A}_{i,j}^o)) \right) + c(\mathcal{A}_{p,q}^h) - c(\mathcal{A}_{p,q}^o) \geq 0$
 - $\Rightarrow \left(\sum_{\substack{1 \leq i < j \leq k \\ (i,j) \neq (p,q)}} (c(\mathcal{A}_{i,j}^h) - c(\mathcal{A}_{i,j}^o)) \right) + c(\mathcal{A}_{p,q}^h) \geq c(\mathcal{A}_{p,q}^o)$

Carrillo-Lipman bound

- $c(\mathcal{A}_{i,j}^o) \geq d(S_i, S_j)$ for any i and j

$$\therefore \left(\sum_{\substack{1 \leq i < j \leq k \\ (i,j) \neq (p,q)}} (c(\mathcal{A}_{i,j}^h) - c(\mathcal{A}_{i,j}^o)) \right) + c(\mathcal{A}_{p,q}^h) \geq c(\mathcal{A}_{p,q}^o)$$

$$\therefore \left(\sum_{\substack{1 \leq i < j \leq k \\ (i,j) \neq (p,q)}} (c(\mathcal{A}_{i,j}^h) - d(S_i, S_j)) \right) + c(\mathcal{A}_{p,q}^h) \geq c(\mathcal{A}_{p,q}^o)$$

- Let $U = \sum_{1 \leq i < j \leq k} c(\mathcal{A}_{i,j}^h)$.
- Let $L = \sum_{1 \leq i < j \leq k} d(S_i, S_j)$.

Carrillo-Lipman bound

$$\therefore \left(\sum_{\substack{1 \leq i < j \leq k \\ (i,j) \neq (p,q)}} (c(\mathcal{A}_{i,j}^h) - d(S_i, S_j)) \right) + c(\mathcal{A}_{p,q}^h) \geq c(\mathcal{A}_{p,q}^o)$$

$$\therefore U - L - (c(\mathcal{A}_{p,q}^h) - d(S_p, S_q)) + c(\mathcal{A}_{p,q}^h) \geq c(\mathcal{A}_{p,q}^o)$$

$$\Rightarrow U - L + d(S_p, S_q) \geq c(\mathcal{A}_{p,q}^o)$$

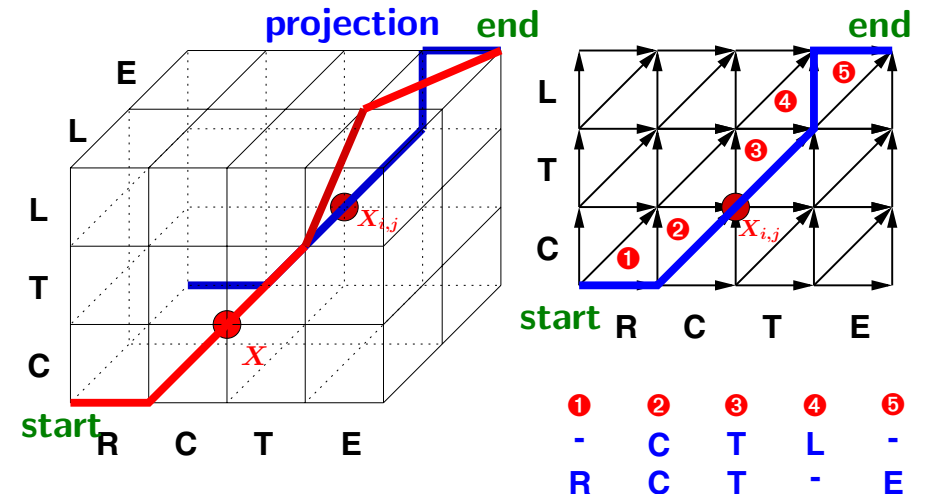
(called Carrillo-Lipman bound, [CL88])

- For an optimal multiple alignment \mathcal{A}^o of S , $c(\mathcal{A}_{i,j}^o) \leq U - L + d(S_i, S_j)$ for all i and j .

Polyhedron

- For an optimal multiple alignment \mathcal{A}^o of S , $c(\mathcal{A}_{i,j}^o) \leq U - L + d(S_i, S_j)$ for all i and j .
- If \mathcal{A} is a candidate of an optimal MSA of S , then $c(\mathcal{A}_{i,j}) \leq U - L + d(S_i, S_j)$ for all i and j .
- **Polyhedron**: consisting of those nodes that are traversed by least one path which corresponds to an alignment \mathcal{A} such that for all i and j , $c(\mathcal{A}_{i,j}) \leq U - L + d(S_i, S_j)$
- How to determine if a node is in the polyhedron?

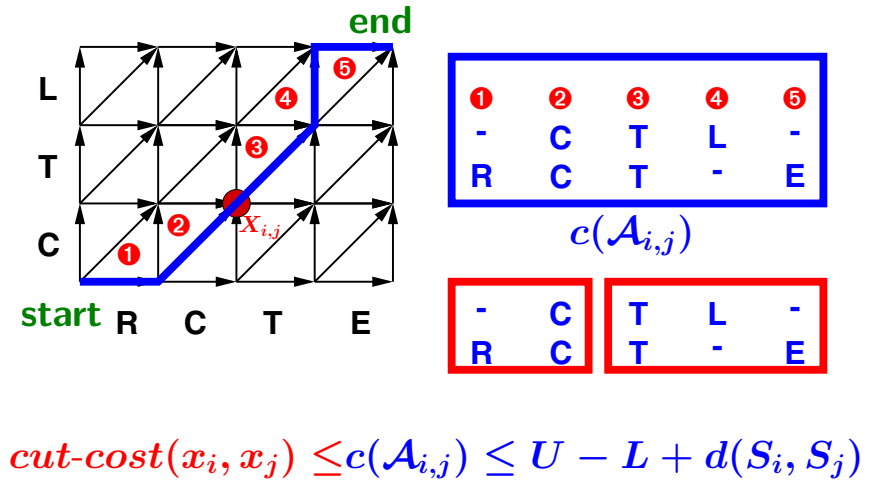
Determine a node of polyhedron



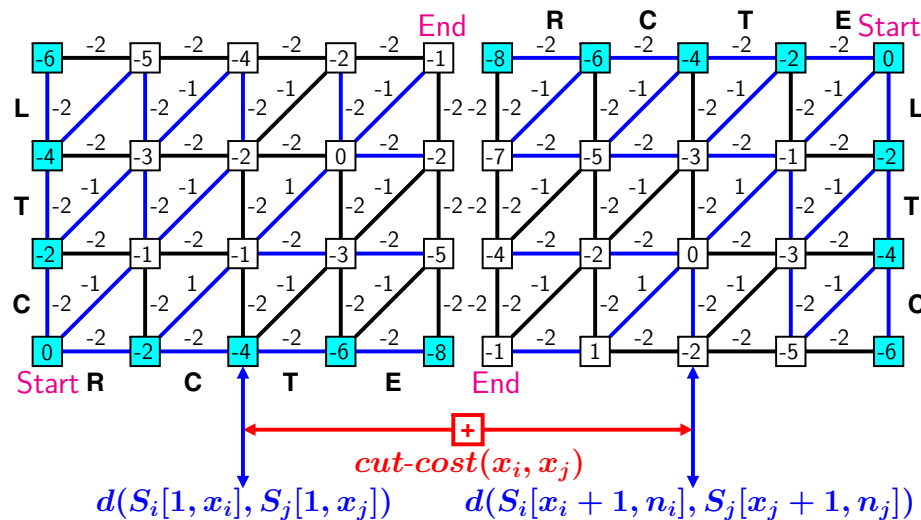
Determine a node of polyhedron

- Let \mathcal{A} be an alignment whose corresponding path P is in the polyhedron.
 - $c(\mathcal{A}_{i,j}) \leq U - L + d(S_i, S_j)$ for any i, j
- Let $X = (x_1, x_2, \dots, x_k)$ be a node in P .
- Let $X_{i,j} = (x_i, x_j)$ be the projected node of X in the projection plane w.r.t. S_i and S_j .
- For $S = s_1 \dots s_n$, let $S[i, j] = s_i \dots s_j$.
- Let $cut-cost(x_i, x_j) = d(S_i[1, x_i], S_j[1, x_j]) + d(S_i[x_i + 1, n_i], S_j[x_j + 1, n_j])$.

Determine a node of polyhedron



Computation of $cut-cost(x_i, x_j)$



Branch and bound method

- Branch (searching) method:** the lexicographic order of the indices of nodes
- Bound rules:** eliminate the currently visiting vertex $X = (x_1, x_2, \dots, x_k)$ if the following condition is satisfied:

Carrillo-Lipman pruning: violate "for all x_i, x_j , $cut-cost(x_i, x_j) \leq U - L + d(S_i, S_j)$ "

- $U = \sum_{1 \leq i < j \leq k} c(\mathcal{A}_{i,j}^h)$
- $L = \sum_{1 \leq i < j \leq k} d(S_i, S_j)$

Sum-of-pairs MSA problem

- NP-complete: by Wang and Jiang [WJ94]; Bonizzoni and Vedova [BV01]
- Branch and bound methods: by Carrillo and Lipman [CL88] (improved by [LAK89, GKS95])
- **Approximate methods:** (k sequences)
 - $(2 - \frac{2}{k})$ -approximation, by Gusfield [Gus93]
 - $(2 - \frac{3}{k})$ -approximation, by Pevzner [Pev92]
 - $(2 - \frac{l}{k})$ -approximation, where $l < k$ by Bafna, Lawler and Pevzner [BLP97]
- Heuristic methods: Progressive algorithms

2-Approximation algorithm ①

- Define the scoring function as follows:

$$\sigma(x, y) = \begin{cases} 0, & \text{if } x = y, \\ 1, & \text{otherwise} \end{cases}$$

- $D(S_i, S_j)$: globally aligned distance between S_i and S_j
- Consider the set S of four sequences:

$S_1 = \text{ATGCTC}$

$S_2 = \text{AGAGC}$

$S_3 = \text{TTCTG}$

$S_4 = \text{ATTGCATGC}$

2-Approximation algorithm ②

$S_1 = \text{ATGCTC}$

$S_2 = \text{AGAGC}$

$S_2 = \text{A-GAGC}$

$S_3 = \text{TTCTG}$

$$D(S_1, S_2) = 3$$

$$D(S_2, S_3) = 5$$

$S_1 = \text{ATGCTC}$

$S_2 = \text{A--G-A-GC}$

$S_3 = \text{TT-CTG}$

$S_4 = \text{ATTGCATGC}$

$$D(S_1, S_3) = 3$$

$$D(S_2, S_4) = 4$$

$S_1 = \text{AT-GC-T-C}$

$S_3 = \text{-TT-C-TG-}$

$S_4 = \text{ATTGCATGC}$

$S_4 = \text{ATTGCATGC}$

$$D(S_1, S_4) = 3$$

$$D(S_3, S_4) = 4$$

2-Approximation algorithm ③

- Find the **center sequence** S_i of S which **minimizes** $\sum_{X \in S \setminus \{S_i\}} D(S_i, X)$.

$$D(S_1, S_2) + D(S_1, S_3) + D(S_1, S_4) = 9$$

$$D(S_2, S_1) + D(S_2, S_3) + D(S_2, S_4) = 12$$

$$D(S_3, S_1) + D(S_3, S_2) + D(S_3, S_4) = 12$$

$$D(S_4, S_1) + D(S_4, S_2) + D(S_4, S_3) = 11$$

Hence, S_1 is the center in this example

2-Approximation algorithm ④

1. Align S_2 with S_1 :
 $S_1 = \text{ATGCTC}$
 $S_2 = \text{A-GAGC}$

2. Add S_3 by aligning S_3 with S_1 :
 $S_1 = \text{ATGCTC}$
 $S_2 = \text{A-GAGC}$
 $S_3 = \text{-TTCTG}$

3. Add S_4 by aligning S_4 with S_1 :
 $S_3 = \text{-T-T-C-T-G}$
 $S_2 = \text{A--GA-G-C}$
 $S_1 = \text{AT-GC-T-C}$
 $S_4 = \text{ATTGCATGC}$

By C.L. Lu

Multiple Sequence Alignment p.33

2-Approximation algorithm ⑤

- $d(S_i, S_j)$ ($d^*(S_i, S_j)$): the distance between S_i and S_j induced by this approximation (an optimal) algorithm
- $d(S_i, S_j) + d(S_i, S_k) \geq d(S_j, S_k)$ (triangle inequality)
- $App = \frac{1}{2} \sum_{i=1}^k \sum_{\substack{j=1 \\ j \neq i}}^k d(S_i, S_j)$
- $Opt = \frac{1}{2} \sum_{i=1}^k \sum_{\substack{j=1 \\ j \neq i}}^k d^*(S_i, S_j)$

By C.L. Lu

Multiple Sequence Alignment p.34

2-Approximation algorithm ⑥

- Claim that $App \leq 2Opt$
- Since triangle inequality property of $d(S_i, S_j)$ and $d(S_1, S_i) = d(S_i, S_1)$, we have

$$\begin{aligned} App &= \frac{1}{2} \sum_{i=1}^k \sum_{\substack{j=1 \\ j \neq i}}^k d(S_i, S_j) \\ &\leq \frac{1}{2} \sum_{i=1}^k \sum_{\substack{j=1 \\ j \neq i}}^k (d(S_i, S_1) + d(S_1, S_j)) \\ &= (k-1) \sum_{i=2}^k d(S_1, S_i) \end{aligned}$$

By C.L. Lu

Multiple Sequence Alignment p.35

2-Approximation algorithm ⑦

- Since $d(S_1, S_i) = D(S_1, S_i)$ for all i , we have

$$App \leq (k-1) \sum_{i=2}^k D(S_1, S_i)$$

- Note that $Opt = \frac{1}{2} \sum_{i=1}^k \sum_{\substack{j=1 \\ j \neq i}}^k d^*(S_i, S_j)$
- Since $D(S_i, S_j) \leq d^*(S_i, S_j)$, we have

$$Opt = \frac{1}{2} \sum_{i=1}^k \sum_{\substack{j=1 \\ j \neq i}}^k d^*(S_i, S_j) \geq \frac{1}{2} \sum_{i=1}^k \sum_{\substack{j=1 \\ j \neq i}}^k D(S_i, S_j)$$

By C.L. Lu

Multiple Sequence Alignment p.36

2-Approximation algorithm

⑧

- Since S_1 is the center, we have

$$\begin{aligned} Opt &\geq \frac{1}{2} \sum_{i=1}^k \sum_{\substack{j=1 \\ j \neq i}}^k D(S_i, S_j) \\ &\geq \frac{1}{2} \sum_{i=1}^k \sum_{j=2}^k D(S_1, S_j) = \frac{k}{2} \sum_{j=2}^k D(S_1, S_j) \end{aligned}$$

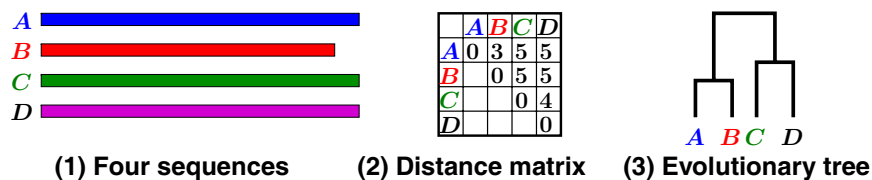
$$App \leq \frac{2(k-1)}{k} \Leftrightarrow App \leq (2 - \frac{2}{k}) Opt < 2 \cdot Opt$$

Sum-of-pairs MSA problem

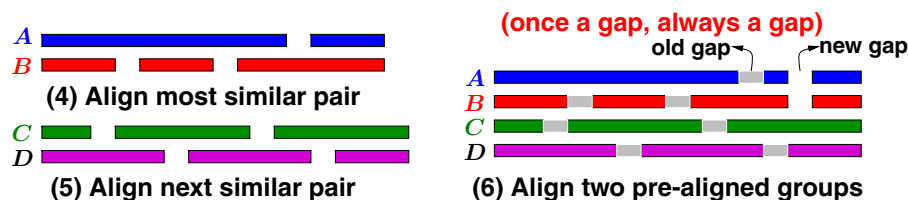
- NP-complete: by Wang and Jiang [WJ94]; Bonizzoni and Vedova [BV01]
- Branch and bound methods: by Carrillo and Lipman [CL88] (improved by [LAK89, GKS95])
- Approximate methods: (k sequences)
 - $(2 - \frac{2}{k})$ -approximation, by Gusfield [Gus93]
 - $(2 - \frac{3}{k})$ -approximation, by Pevzner [Pev92]
 - $(2 - \frac{l}{k})$ -approximation, where $l < k$ by Bafna, Lawler and Pevzner [BLP97]
- Heuristic methods: Progressive algorithms

Progressive MSA

- Construct the guide (evolutionary) tree



- Progressively align sequences by the tree



Progressive MSA

- Compute the distance matrix by aligning all pairs of sequences (using dynamic programming algorithm)
- Compute the guide tree from the distance matrix:
 - PILEUP of GCG: UPGMA (Unweighted Pair-Group Method using Arithmetic mean)
 - CLUSTAL W: NJ (Neighbor-Joining)
 - YAMA-MST of Tang's lab: Kruskal MST
- Progressively align the sequences according to the branching order in the guide tree (once a gap, always a gap)

1. Computation of distance matrix

S_1 : RCTLEE
 S_2 : RCLEE
 S_3 : CTLEE
 S_4 : CTEE

	S_1	S_2	S_3	S_4
S_1	0	1	1	2
S_2		0	2	1
S_3			0	1
S_4				0

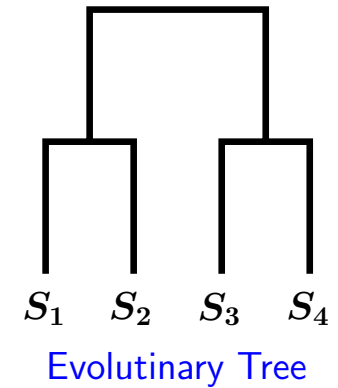
4 Sequences →

Distance Matrix

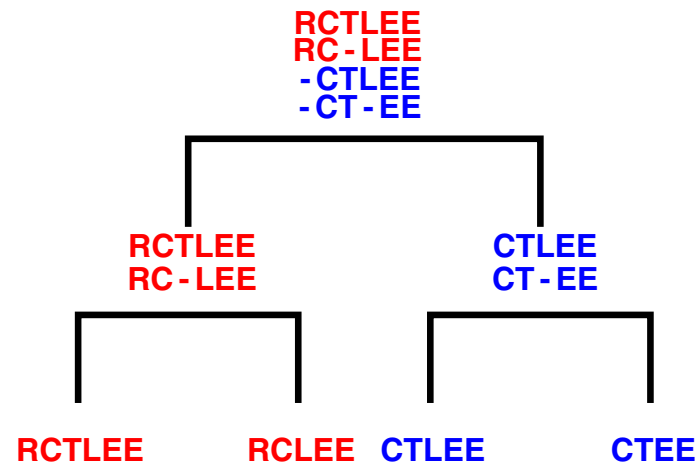
2. Construction of guide tree

	S_1	S_2	S_3	S_4
S_1	0	1	1	2
S_2		0	2	1
S_3			0	1
S_4				0

Distance Matrix →



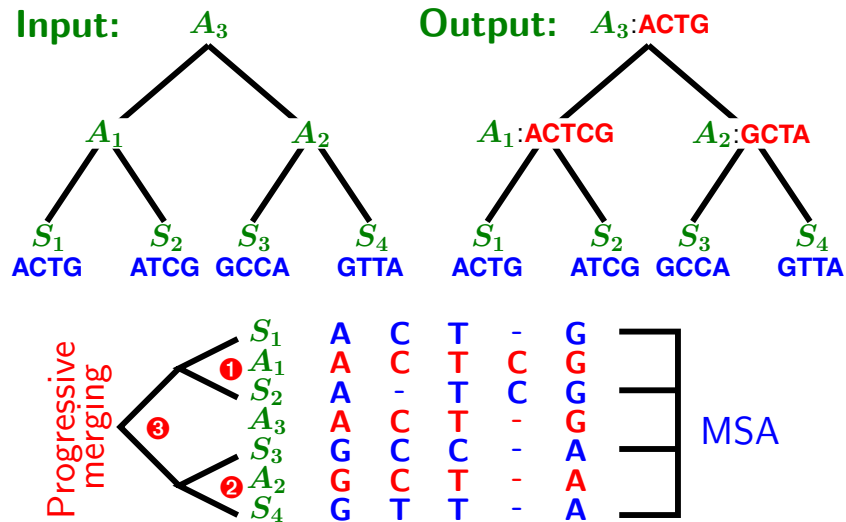
3. Progressive alignment



Tree alignment problem

- **Input:** k sequences and a rooted phylogenetic tree containing k leaves, each of which is labeled with a unique given sequence
- **Output:** Construct a sequence for each internal node of the tree such that the total cost of its edges is minimized, where the cost of an edge is the edit distance between two sequences associated with both ends of the edge
- The given tree represents the evolutionary history for the leaf sequences of extant species, and its internal nodes represent the extinct ancestral species whose sequences are to be determined.

Tree alignment: example



By C.L. Lu

Multiple Sequence Alignment p.45

Tree alignment: complexity

- **NP-complete:** Wang and Jiang, [WJ94]
- **2-approximation:** Wang, Jiang and Lawler, [WJL96]
- **PTAS:** Wang, Jiang and Lawler, [WJL96] (further improved by [WG97, WJG00])
- **Polynomial Time Approximation Scheme:** a family of algorithms $\{\mathcal{A}_\epsilon : \epsilon > 0\}$ each of which takes as **input** both an instance I and an error bound ϵ and in polynomial time, it outputs an approximate solution $\mathcal{A}(I)$ whose **performance ratio** is

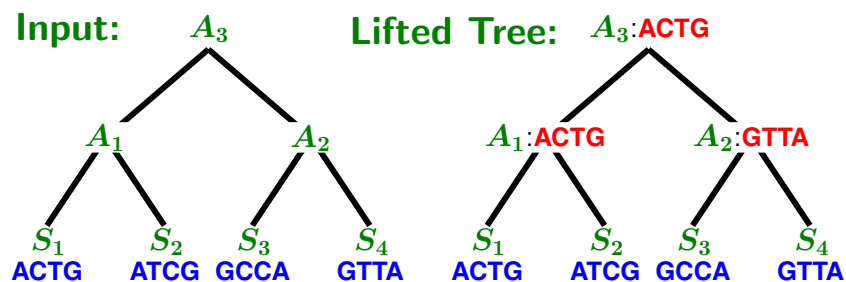
$$R_{\mathcal{A}}(I, \epsilon) = \max\left\{\frac{\mathcal{A}(I)}{OPT(I)}, \frac{OPT(I)}{\mathcal{A}(I)}\right\} \leq 1 + \epsilon$$

By C.L. Lu

Multiple Sequence Alignment p.46

Tree alignment: 2-approximation

- **Lifted tree:** an evolutionary tree in which the label of each internal node equals to the label of some one of its children



By C.L. Lu

Multiple Sequence Alignment p.47

Tree alignment: 2-approximation

- **Theorem:** There exists a lifted tree with cost at most $2(1 - \frac{1}{k})c(T_{opt})$.
 - k : the number of leaves of the given tree T
 - T_{opt} : an optimal tree with cost $c(T_{opt})$
- **Closest descendant leaf $l(v)$ of $v \in T_{opt}$:** a descendant leaf of v s.t. the path from v to $l(v)$ is the shortest among all descendant leaves of v
- T_l : a lifted tree of T obtained by assigning the sequence $sl(v)$ of $l(v)$ to each internal node v
- **Lemma:** $c(T_l) \leq 2(1 - \frac{1}{k})c(T_{opt})$

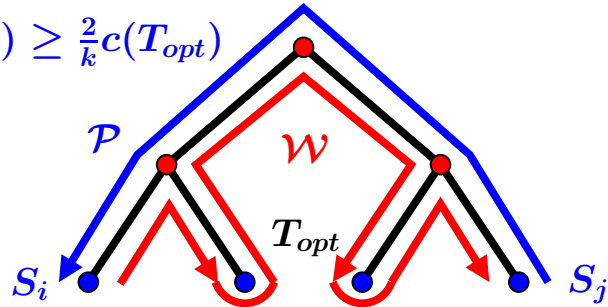
By C.L. Lu

Multiple Sequence Alignment p.48

$$c(T_l) \leq 2\left(1 - \frac{1}{k}\right)c(T_{opt}) \quad \textcircled{1}$$

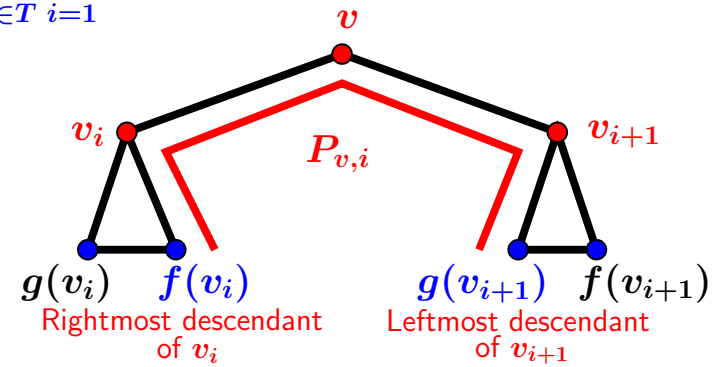
- \mathcal{P} : the longest path among all paths between any two leaves S_i and S_j in T_{opt}
 - \mathcal{W} : the walk obtained by removing \mathcal{P} from an Euler tour of T_{opt}
- $\therefore c(\mathcal{W}) \leq 2\left(1 - \frac{1}{k}\right)c(T_{opt})$

$$\therefore c(\mathcal{P}) \geq \frac{2}{k}c(T_{opt})$$



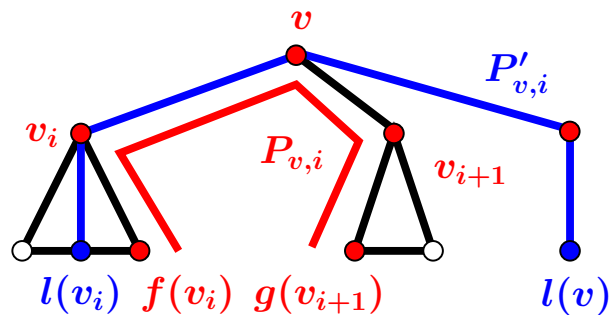
$$c(T_l) \leq 2\left(1 - \frac{1}{k}\right)c(T_{opt}) \quad \textcircled{2}$$

- $P_{v,i}$: path $f(v_i) \rightarrow v_i \rightarrow v \rightarrow v_{i+1} \rightarrow g(v_{i+1})$
- $\sum_{v \in T} \sum_{i=1}^{d-1} c(P_{v,i}) = c(\mathcal{W})$ (d : # of v 's children)



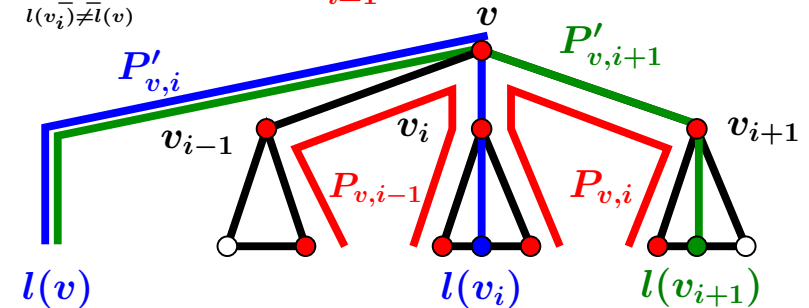
$$c(T_l) \leq 2\left(1 - \frac{1}{k}\right)c(T_{opt}) \quad \textcircled{3}$$

- $P'_{v,i}$: path $l(v) \rightarrow v \rightarrow v_i \rightarrow l(v_i)$
- $c(P'_{v,i}) \leq c(P_{v,i})$ and $c(P'_{v,i+1}) \leq c(P_{v,i})$



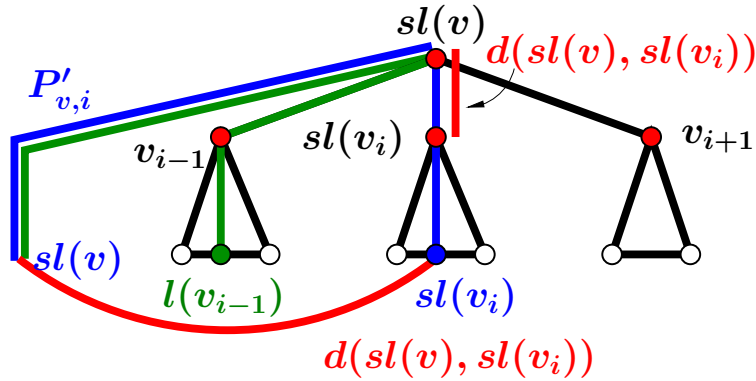
$$c(T_l) \leq 2\left(1 - \frac{1}{k}\right)c(T_{opt}) \quad \textcircled{4}$$

- $c(P'_{v,i}) \leq c(P_{v,i-1})$ and $c(P'_{v,i+1}) \leq c(P_{v,i})$
- $\sum_{\substack{1 \leq i \leq d \\ l(v_i) \neq l(v)}} c(P'_{v,i}) \leq \sum_{i=1}^{d-1} c(P_{v,i})$



$$c(T_l) \leq 2\left(1 - \frac{1}{k}\right)c(T_{opt}) \quad \textcircled{5}$$

- The cost of edge $(v, v_i) = d(sl(v), sl(v_i))$, where $d(sl(v), sl(v_i)) = 0$ if $l(v) = l(v_i)$
- By triangle inequality, $d(sl(v), sl(v_i)) \leq c(P'_{v,i})$



By C.L. Lu

Multiple Sequence Alignment p.53

$$c(T_l) \leq 2\left(1 - \frac{1}{k}\right)c(T_{opt}) \quad \textcircled{6}$$

- The cost of T_l is as follows:

$$\begin{aligned} c(T_l) &= \sum_{v \in T} \sum_{i=1}^d d(sl(v), sl(v_i)) \\ &\leq \sum_{v \in T} \sum_{\substack{1 \leq i \leq d \\ l(v_i) \neq l(v)}} c(P'_{v,i}) \\ &\leq \sum_{v \in T} \sum_{i=1}^{d-1} c(P_{v,i}) \\ &= c(\mathcal{W}) \leq 2\left(1 - \frac{1}{k}\right)c(T_{opt}) \end{aligned}$$

By C.L. Lu

Multiple Sequence Alignment p.54

Tree alignment: 2-approximation

- Computing T_l is hard since it is derived from T_{opt} and the computation of T_{opt} is NP-complete.
- T^* : an optimal one among all the lifted trees of T .
- Clearly, $c(T^*) \leq c(T_l) \leq 2\left(1 - \frac{1}{k}\right)c(T_{opt})$
- How to compute T^* ? (By dynamic programming)
 - T_v : the induced subtree of T rooted at v
 - $S(v)$: the set of leaf sequences of T_v
 - $D[v, S_i]$: the cost of an optimal lifted tree of T_v with v being assigned $S_i \in S(v)$
 - $c(T^*) = \min\{D[root, S_i] : 1 \leq i \leq k\}$

By C.L. Lu

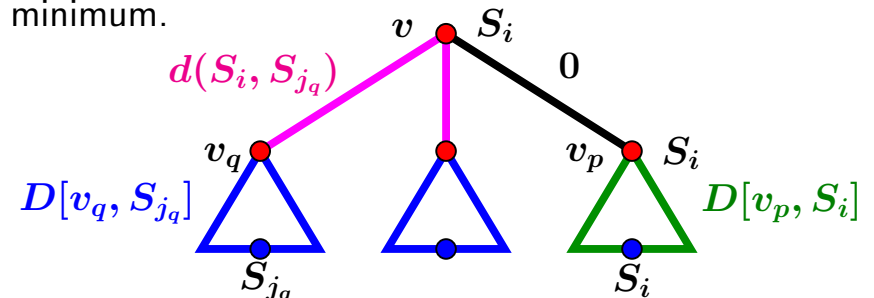
Multiple Sequence Alignment p.55

How to compute $D[v, S_i]$?

Let v be an internal node of T with v_1, \dots, v_d children, and let $S_i \in S(v_p)$, $1 \leq p \leq d$. Then

$$D[v, S_i] = D[v_p, S_i] + \sum_{\substack{1 \leq q \leq d \\ q \neq p}} (D[v_q, S_{j_q}] + d(S_i, S_{j_q}))$$

where S_{j_q} is in $S(v_q)$ s.t. $D[v_q, S_{j_q}] + d(S_i, S_{j_q})$ is minimum.



By C.L. Lu

Multiple Sequence Alignment p.56

- [LAK89] D. J. Lipman, S. F. Altschul, and J. D. Kececioglu. A tool for multiple sequence alignment. *Proceedings of the National Academy of Sciences USA*, 86:4412–4415, 1989.
- [Pev92] P. A. Pevzner. Multiple alignment, communication cost, and graph matching. *SIAM Journal on Applied Mathematics*, 52(6):1763–1779, 1992.
- [WG97] L. Wang and D. Gusfield. Improved approximation algorithms for tree alignment. *Journal of Algorithms*, 25:255–273, 1997.
- [WJ94] L. Wang and T. Jiang. On the complexity of multiple sequence alignment. *Journal of Computational Biology*, 1:337–348, 1994.
- [WJG00] L. Wang, T. Jiang, and D. Gusfield. A more efficient approximation scheme for tree alignment. *SIAM Journal on Computing*, 30(1):283–299, 2000.
- [WJL96] L. Wang, T. Jiang, and E. L. Lawler. Approximation algorithms for tree alignment with a given phylogeny. *Algorithmica*, 16(3):302–315, 1996.

56-2

References

- [BLP97] V. Batna, E. L. Lawler, and P. A. Pevzner. Approximation algorithms for multiple sequence alignment. *Theoretical Computer Science*, 182(1-2):233–244, 1997.
- [BV01] P. Bonizzoni and G. D. Vedova. The complexity of multiple sequence alignment with SP-score that is a metric. *Theoretical Computer Science*, 259(1–2):63–79, 2001.
- [CL88] H. Carrillo and D. Lipman. The multiple sequence alignment problem in biology. *SIAM Journal on Applied Mathematics*, 48:1073–1082, 1988.
- [GKS95] S. K. Gupta, J. Kececioglu, and A. A. Schäffer. Improving the practical space and time efficiency of the shortest-paths approach to sum-of-pairs multiple sequence alignment. *Journal of Computational Biology*, 2:459–472, 1995.
- [Gus93] D. Gusfield. Efficient methods for multiple sequence alignment with guaranteed error bounds. *Bulletin of Mathematical Biology*, 55(1):141–154, 1993.